

Large Concept Models:

Language Modeling in a Sentence Representation Space

The LCM team, Loïc Barrault*, Paul-Ambroise Duquenne*, Maha Elbayad*, Artyom Kozhevnikov*, Belen Alastruey†, Pierre Andrews†, Mariano Coria†, Guillaume Couairon^{††}, Marta R. Costa-jussà†, David Dale†, Hady Elsahar†, Kevin Heffernan†, João Maria Janeiro†, Tuan Tran†, Christophe Ropers†, Eduardo Sánchez†, Robin San Roman†, Alexandre Mourachko‡, Safiyyah Saleem‡, Holger Schwenk‡

FAIR at Meta

*Core contributors, alphabetical order, †Contributors to data preparation, LCM extensions and evaluation, alphabetical order, ‡Research and project management, alphabetical order, ^{††}Initial work while at FAIR at Meta, new affiliation: INRIA, France

LLMs have revolutionized the field of artificial intelligence and have emerged as the de-facto tool for many tasks. The current established technology of LLMs is to process input and generate output at the token level. This is in sharp contrast to humans who operate at multiple levels of abstraction, well beyond single words, to analyze information and to generate creative content. In this paper, we present an attempt at an architecture which operates on an explicit higher-level semantic representation, which we name a “*concept*”. Concepts are language- and modality-agnostic and represent a higher level idea or action in a flow. Hence, we build a “*LARGE CONCEPT MODEL*”. In this study, as proof of feasibility, we assume that a concept corresponds to a sentence, and use an existing sentence embedding space, SONAR, which supports up to 200 languages in both text and speech modalities.

The *LARGE CONCEPT MODEL* is trained to perform autoregressive sentence prediction in an embedding space. We explore multiple approaches, namely MSE regression, variants of diffusion-based generation, and models operating in a quantized SONAR space. These explorations are performed using 1.6B parameter models and training data in the order of 1.3T tokens. We then scale one architecture to a model size of 7B parameters and training data of about 2.7T tokens. We perform an experimental evaluation on several generative tasks, namely summarization and a new task of summary expansion. Finally, we show that our model exhibits impressive zero-shot generalization performance to many languages, outperforming existing LLMs of the same size. The training code of our models is freely available.^a

Date: December 12, 2024

Correspondence: Holger Schwenk at schwenk@meta.com

^ahttps://github.com/facebookresearch/large_concept_model



1 Introduction

Large Language models (LLMs) are dominating current research in natural language processing, and with their recent extension to more modalities, namely images, video and speech, they seem to be considered as the de-facto technique to follow to approach human intelligence. LLMs achieve indeed impressive performance on a large variety of tasks, such as providing detailed answers for general knowledge questions, helping in performing long document analysis, or drafting different types of messages, and writing or debugging code. Building an LLM from scratch requires access to

enormous computational resources to process ever larger amounts of data and train models, the size of which now exceeds four hundred billion parameters. Knowledge acquisition in LLMs is heavily data-driven and extending them to more languages or modalities usually requires injecting additional (synthetic) data to cover them.

The landscape of available LLMs can be structured into open models such as LLAMA (The Llama3 team, 2024), MISTRAL (Jiang et al., 2024), BLOOM (BigScience Workshop, 2023) or FALCON (Almazrouei et al., 2023), on the one hand, and closed models such as GEMINI (Gemini Team Google, 2024), GPT (OpenAI, 2024) or CLAUDE (Anthropic, 2024), on the other. It is striking that all these models are based on the same underlying architecture: a transformer-based, decoder-only language model, which is pretrained to predict the next token, given a long context of preceding tokens. Despite the undeniable success of LLMs and continued progress, all current LLMs miss a crucial characteristic of human intelligence: explicit reasoning and planning at multiple levels of abstraction. The human brain does not operate at the word level only. We usually have a top-down process to solve a complex task or compose a long document: we first plan at a higher level the overall structure, and then step-by-step, add details at lower levels of abstraction. One may argue that LLMs are implicitly learning a hierarchical representation, but we stipulate that models with an explicit hierarchical architecture are better suited to create coherent long-form output.

Imagine a researcher giving a fifteen-minute talk. In such a situation, researchers do not usually prepare detailed speeches by writing out every single word they will pronounce. Instead, they outline a flow of higher-level ideas they want to communicate. Should they give the same talk multiple times, the actual words being spoken may differ, the talk could even be given in different languages, but the flow of higher-level abstract ideas will remain the same. Similarly, when writing a research paper or essay on a specific topic, humans usually start by preparing an outline that structures the whole document into sections, which they then refine iteratively. Humans also detect and remember dependencies between the different parts of a longer document at an abstract level. If we expand on our previous research writing example, keeping track of dependencies means that we need to provide results for each of the experiment mentioned in the introduction. Finally, when processing and analyzing information, humans rarely consider every single word in a large document. Instead, we use a hierarchical approach: we remember which part of a long document we should search to find a specific piece of information.

To the best of our knowledge, this explicit hierarchical structure of information processing and generation, at an abstract level, independent of any instantiation in a particular language or modality, cannot be found in any of the current LLMs. In this work, we present a new approach which moves away from processing at the token level and closer to (hierarchical) reasoning in an abstract embedding space. This abstract embedding space is designed to be independent of the language or modality in which the content is expressed; in other words, we aim to model the underlying reasoning process at a purely semantic level, not its instantiation in a specific language. In order to verify our approach, we limit our study to two levels of abstraction: subword tokens and *concepts*. We define a *concept* as an abstract atomic idea. In practice, a concept would often correspond to a sentence in a text document, or an equivalent speech utterance. We posit that a sentence is an appropriate unit to achieve language independence, in opposition to single words. This is in sharp contrast to current LLMs techniques which are heavily English centric and token based.

Our fundamental idea could be based on any fixed-size sentence embedding space for which an encoder and decoder are available. In particular, we could aim to train a new embedding space specifically optimized to our reasoning architecture. In this work, we chose an existing and freely available sentence embedding, named SONAR (Duquenne et al., 2023b). SONAR supports text

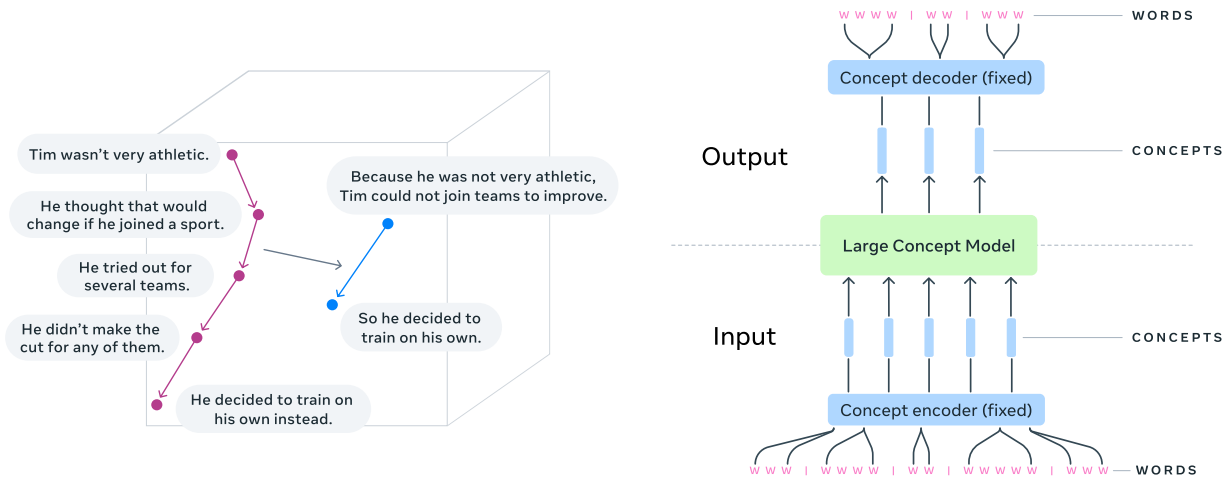


Figure 1 - Left: visualization of reasoning in an embedding space of concepts (task of summarization). Right: fundamental architecture of an LARGE CONCEPT MODEL (LCM).
 *: concept encoder and decoder are frozen.

input and output in 200 languages, speech input in 76 languages, and speech output in English. We discuss the constraints and impact of this choice in [Section 2.1](#), and share some ideas on alternative embedding spaces in [Section 6](#).

[Figure 1](#)-left visualizes reasoning in an embedding space with the example of a summarization task, which is materialized by a function on the embedding space, mapping five concept representations into two. [Figure 1](#)-right summarizes the overall architecture and processing flow. The input is first segmented into sentences, and each one is encoded with SONAR to achieve a sequence of concepts, *i.e.*, sentence embeddings. This sequence of concepts is then processed by a **Large Concept Model (LCM)** to generate at the output a new sequence of concepts. Finally, the generated concepts are decoded by SONAR into a sequence of subwords. The encoder and decoder are fixed and are not trained. It is important to highlight that the unchanged sequence of concepts at the output of the LCM can be decoded into other languages or modalities without performing again the whole reasoning process. In the same spirit, a particular reasoning operation such as summarization can be performed in a zero-shot setting on input in any language or modality, since it solely operates on concepts. To summarize, the LCM neither has information on the input language or modality nor generates output in a particular language or modality. We explore multiple architectures to train the LCM, in particular several variants of diffusion. Finally, we envision an additional level of abstraction beyond concepts which could correspond to a short description of a paragraph or small section. In [Section 4.3](#) we report initial ideas on how conditioning and predicting such higher-level representations can improve consistency of output generated by an LCM.

To some extent, the LCM architecture resembles the JEPa approach ([LeCun, 2022](#)) that also aims to predict the representation of the next observation in an embedding space. However, unlike JEPa that places more emphasis on learning a representation space in a self-supervised way, the LCM focuses on accurate prediction in the existing embedding space.

The main characteristics of our generic LARGE CONCEPT MODEL approach are as follows:

- **Reasoning at an abstract language- and modality-agnostic level beyond tokens:**
 - We model the underlying reasoning process, not its instantiation in a particular language.
 - The LCM can be trained, i.e. acquire knowledge, on all languages and modalities at once, promising scalability in an unbiased way.
- **Explicit hierarchical structure:**
 - Better readability of long-form output by a human.
 - Facilitates local interactive edits by a user.
- **Handling of long context and long-form output:**
 - The complexity of a vanilla transformer model increases quadratically with the sequence length. This makes handling of large context windows challenging and several techniques have been developed to alleviate this problem, *e.g.*, sparse attention (Child et al., 2019) or LSH attention (Kitaev et al., 2020). Our LCM operates on sequences which are at least an order of magnitude shorter.¹
- **Unparalleled zero-shot generalization:**
 - Independently of the language or modality the LCM is pre-trained and fine-tuned on, it can be applied to any language and modality supported by the SONAR encoders, without the need of additional data or fine-tuning. We report results for multiple languages in the text modality.
- **Modularity and extensibility:**
 - Unlike multimodal LLMs that can suffer from modality competition (Aghajanyan et al., 2023; Chameleon team, 2024), concept encoders and decoders can be independently developed and optimized without any competition or interference.
 - New languages or modalities can be easily added for an existing system.

The goal of this paper is to provide a proof of concept of this high-level vision of an alternative architecture to current best practice in language modeling. In the next section we present the main design principles of our models and discuss several variants to build and train a LARGE CONCEPT MODEL. We discuss several designs to implement diffusion approaches with concept embeddings and carefully study noise scheduling. This section is completed by a compute complexity comparison with token-based LLMs. Section 3 is dedicated to the analysis of a larger 7B parameter model. We discuss challenges when instruction fine-tuning this model on multiple generative tasks, and provide a comparison with existing LLMs of comparable size. The paper concludes with a discussion of related work, the current limitations and perspectives of our approach.

To foster research in this area, we make our LCM training code² as well as SONAR encoders and decoders³ for up to 200 languages and multiple modalities freely available.

¹We assume an average sentence length of 10–20 tokens.

²https://github.com/facebookresearch/large_concept_model

³<https://github.com/facebookresearch/SONAR>

2 Main Design Principles

In this section, we outline the main design principles of the LCM. We first describe the SONAR embedding space with its encoders and decoders. Then, we discuss details of data preparation, namely sentence segmentation *i.e.*, how we split long documents into sentences. And finally, we describe in details the different versions of LCMS introduced in this work.

2.1 The SONAR embedding space

The motivation of this work is to perform reasoning at a higher conceptual level than tokens. This requires an embedding space which is highly semantic. We chose SONAR (Duquenne et al., 2023b) since it achieves best performance on several semantic similarity metrics like `xsim` or `xsim++` (Chen et al., 2023b), and it was successfully used in large-scale bitext mining for translation (Seamless Communication et al., 2023b).

The SONAR text embedding space was trained as an encoder/decoder architecture, with a fixed-size bottleneck instead of cross-attention (see Figure 2). The criterion combines a machine translation objective for 200 languages into and out of English, denoising auto-encoding and an explicit MSE loss at the embedding bottleneck layer. Once the text embedding space was trained, a teacher-student approach was applied to extend the SONAR space to the speech modality. More details on the architecture and training procedure can be found in Duquenne et al. (2023b), and detailed speech recognition and translation results in the appendix of Seamless Communication et al. (2023a).

Our LCM operates directly on SONAR concepts embeddings, hence, it can perform reasoning on all supported languages and modalities. Table 1 compares the language coverage of several other LLMs. The LCM supports substantially more languages than other models, in particular many low-resource languages. In addition to the text modality, SONAR supports 76 languages for speech input and speech output in English. We have also developed an experimental encoder for American Sign language (ASL). All these encoders and decoders are freely available.⁴ Exact listings of the supported languages can be found in the SONAR GitHub repository.

⁴<https://github.com/facebookresearch/SONAR>

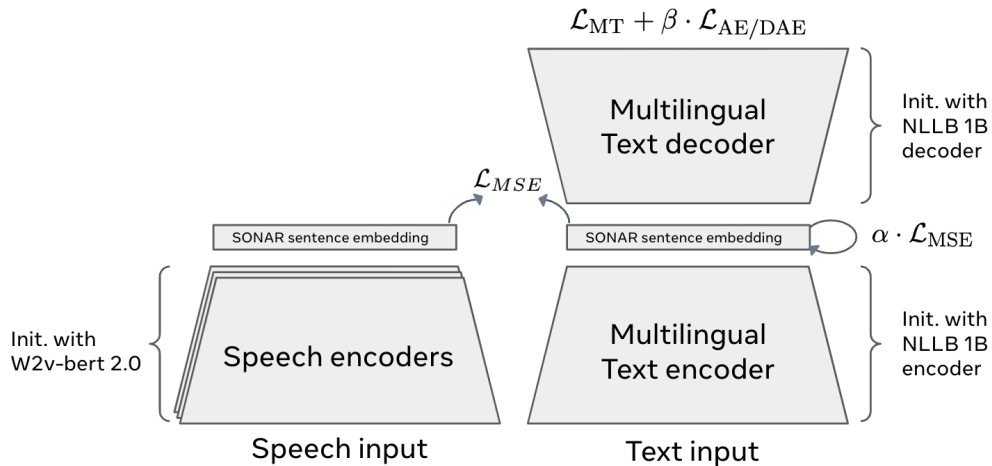


Figure 2 - Encoder/decoder bottleneck architecture to train the SONAR text embeddings (right part of figure). Teacher-student approach to extend SONAR to the speech modality (left part).

Model	Text		Speech		Image		Video	
	Input	Output	Input	Output	Input	Output	Input	Output
GEMINI	47	47	62	✓	✓	✓	✓	✗
GPT	85	85	✓	✓	✓	✓	?	✗
CLAUDE	37	37	✓	✓	✓	✓	✗	✗
BLOOM	46	46	✗	✗	✓	✓	✗	✗
LLAMA 3-400B	8	8	34	✗	✓	✓	✗	✗
LCM-SONAR	200	200	76	1	✗	✗	(ASL)	✗

Table 1 - Comparison of language and modality coverage for several LLMs and our LCM operating on the SONAR embedding space. SONAR has an experimental support for American Sign Language (ASL) which is not used in this paper.

2.2 Data preparation

To train and evaluate the LCM, we need to convert raw text datasets into a sequence of SONAR embeddings, each one corresponding to a sentence. Dealing with large text corpora presents several practical limitations. First, the precise segmentation of a text into sentences can be challenging due to the presence of errors, specific formatting issues or any other sources of noise. This requires us to apply robust automatic text segmentation techniques. Second, some sentences (even well formed) can be very long and complex, which might negatively impact the quality of the encoded SONAR embeddings. This is particularly prevalent for texts in the scientific domain. In the following, we discuss strategies for sentence segmentation and how they affect the SONAR encoding.

Sentence segmentation analysis We have identified two potential sentence segmentation techniques; as we are exploring multilingual data, we focus on sentence segmenters with a large language coverage:

1. Spacy segmenter (SPACY) (Honnibal et al., 2020) is a well established multilingual NLP toolkit that provides a rule-based approach to sentence segmentation. SPACY is thoroughly tested for high-resource languages.
2. Segment any Text (SAT) (Minixhofer et al., 2023; Frohmann et al., 2024) offers a suite of models and adapters that predict sentence boundaries at the token level. SAT is designed to be resilient to perturbations, particularly avoiding the over-reliance on punctuation and capitalization. This is valuable in domains where these conventional markers are often missing. The quality of SAT’s segmentation is however dependent on the choice of an “appropriate” split probability threshold.

We additionally customize both methods by incorporating a maximum sentence length cap in characters. We refer to these extensions by SPACY Capped and SAT Capped. Long sentences are broken down into smaller, logically coherent fragments using a rule-based approach based on punctuation marks for SPACY. For SAT, we leverage the provided splitting probability estimates to identify the next best potential split.

To measure the efficacy of a given segmenter, we evaluate the quality of the reconstructed sentences with AUTOBLEU. It is defined as a BLEU score (Papineni et al., 2002) comparing the decoded text from a SONAR vector after encoding a segment, to the the reference segment. A good segmentation will yield segments that can be encoded and then decoded without loss of signal, and thus score a higher AUTOBLEU.

For this analysis, we sample 10k documents from our pretraining datasets, representing approximately 500k sentences. The documents are processed with each segmenter, the sentences are encoded then decoded and the AUTOBLEU score is calculated. We stratified the results based on the lengths of the original sentences.

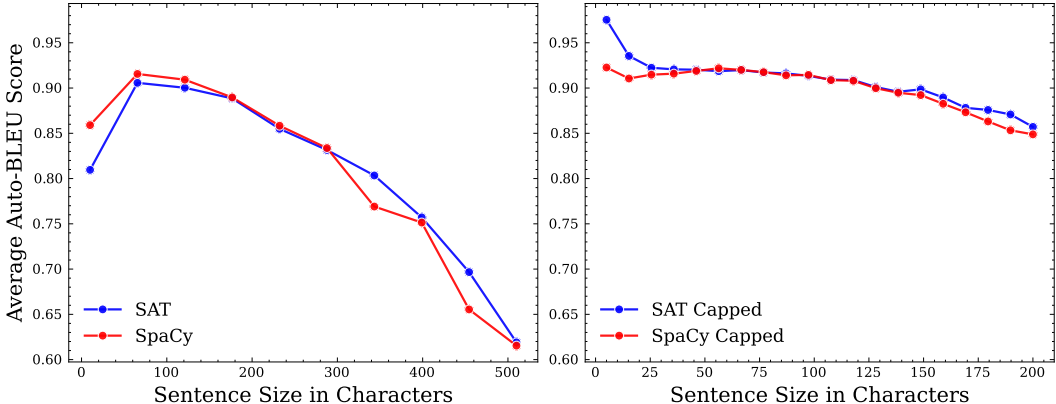


Figure 3 - Segmenters quality. Average Auto-BLEU scores for different sentence segmentation methods depending on sentence length, for both out of the box (left) and capped implementations (right).

As illustrated in Figure 3 and with a capping at 200 characters, the SAT Capped method demonstrates a slight but consistent advantage over SPaCy Capped. Both out-of-the-box segmenters, however, exhibit significant under-performance across all sentence lengths. This lower performance is especially pronounced for sentences exceeding 250 characters, underscoring the limitations of using the segmenters without capping.

Accordingly, we prepare the LCM training data with SAT Capped. We discuss in Appendix A technical and engineering challenges faced when handling large amounts of SONAR embeddings.

2.3 Large Concept Model variants

The design of the LCM is driven by the need to conditionally generate a continuous sentence embedding. This obviously contrasts with how current LLMs work, *i.e.*, estimating a probability distribution over a vocabulary of discrete tokens. A straightforward way of solving the task is to train a transformer model to generate an embedding with the objective of minimizing the MSE loss (see Section 2.3.1). However, a given context may have many plausible, yet semantically different, continuations. The model should thus be able to learn a conditional probability distribution over the continuous embedding of the next sentence.

There is a large body of work in computer vision aiming to learn such conditional probability distributions over continuous data (Dhariwal and Nichol, 2021; Rombach et al., 2021). Models like Dall-E 3 (Betker et al., 2023) or Imagen Video (Ho et al., 2022) use a diffusion process to generate an image or video from a text prompt. Many different real images may satisfy the same input prompt, hence the model has to learn a probability distribution over continuous pixel data. This motivates the exploration of diffusion models for sentence embedding generation. Two variants are presented in Sections 2.3.3 and 2.3.4. Another prevalent take on continuous data generation consists of quantizing said data to ultimately model with discrete units; we explore LCM modeling with quantization in Section 2.3.5.

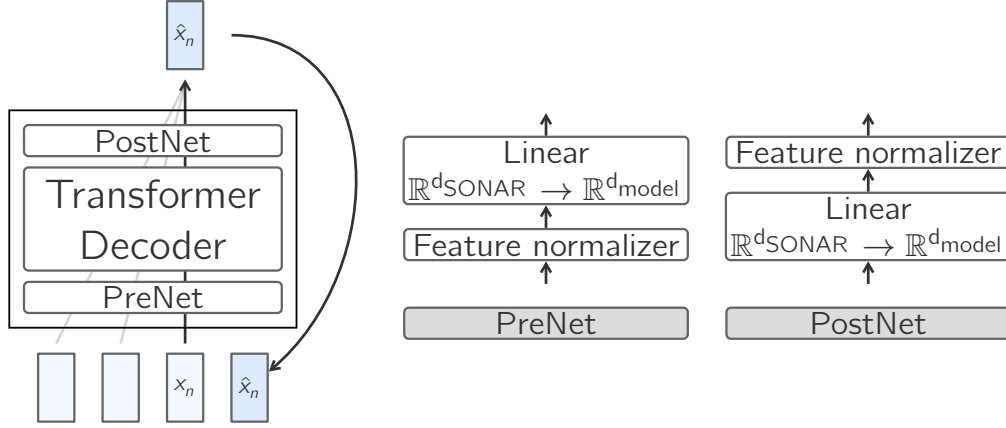


Figure 4 - TheBase-LCM. Illustration of the BASE-LCM. At its core is a standard decoder-only Transformer surrounded with a PreNet and a PostNet.

2.3.1 Base-LCM

Our baseline architecture for next-concept prediction is a standard decoder-only Transformer that transduces a sequence of preceding concepts (read sentence embeddings) into a sequence of future ones. As illustrated in Figure 4, the BASE-LCM is equipped with a “PostNet” and a “PreNet”. The PreNet normalizes the input SONAR embeddings and maps them to the model’s hidden dimension d_{model} .

$$\text{PreNet}(\mathbf{x}) = \text{normalize}(\mathbf{x})\mathbf{W}_{\text{pre}}^t + \mathbf{b}_{\text{pre}}, \quad (1)$$

$$\text{PostNet}(\mathbf{x}) = \text{denormalize}(\mathbf{x}\mathbf{W}_{\text{post}}^t + \mathbf{b}_{\text{post}}), \quad (2)$$

$$(3)$$

where $\mathbf{W}_{\text{post}} \in \mathbb{R}^{d_{\text{SONAR}} \times d_{\text{model}}}$, $\mathbf{b}_{\text{post}} \in \mathbb{R}^{d_{\text{SONAR}}}$, $\mathbf{W}_{\text{pre}} \in \mathbb{R}^{d_{\text{model}} \times d_{\text{SONAR}}}$ and $\mathbf{b}_{\text{pre}} \in \mathbb{R}^{d_{\text{model}}}$.

In order to learn the maps “normalize” and its inverse “denormalize” we fit a robust scaler to a set of randomly sampled SONAR vectors from different corpora and domains of text data. This scaler removes the median statistics and scales the data according to the interquartile range (IQR).

$$\text{normalize}(\mathbf{x}) = \frac{\mathbf{x} - \boldsymbol{\mu}}{\boldsymbol{\sigma}}, \quad \text{denormalize}(\mathbf{x}) = \boldsymbol{\mu} + \boldsymbol{\sigma}\mathbf{x}. \quad (4)$$

The BASE-LCM is trained on the semi-supervised task of next concept prediction, that is, the model predicts the next concept $\hat{\mathbf{x}}_n$ and its parameters $\boldsymbol{\theta}$ are optimized to regress the ground truth next concept (\mathbf{x}_n).

$$\hat{\mathbf{x}}_n = f(\mathbf{x}_{<n}; \boldsymbol{\theta}), \quad \text{MSE}(\hat{\mathbf{x}}_n, \mathbf{x}_n) = \|\hat{\mathbf{x}}_n - \mathbf{x}_n\|^2. \quad (5)$$

Given a data distribution q of documents (sequences of concepts), the training loss is evaluated as:

$$\mathcal{L}_{\text{BASE-LCM}}(\boldsymbol{\theta}) = \mathbb{E}_{\mathbf{x} \sim q} \left[\sum_{n=1}^{|\mathbf{x}|} \text{MSE}(f(\mathbf{x}_{<n}; \boldsymbol{\theta}), \mathbf{x}_n) \right]. \quad (6)$$

In order to enable the generation of variable length documents at inference time, we suffix training documents with the sentence “End of text.”. Similar to any sentence in the document, this special

suffix will be encoded with SONAR. This means that $\mathbf{x}_{|\mathbf{x}|} = \overrightarrow{\text{eot}} := \text{encode}(\text{"End of text."})$. During inference, we implement two main early stopping mechanisms: the first one measures the similarity of the generated embedding $\hat{\mathbf{x}}_n$ to $\overrightarrow{\text{eot}}$ and stops if the cosine similarity exceeds a threshold s_{eot} . The second mechanism compares the newly generated embedding $\hat{\mathbf{x}}_n$ to the previous generation $\hat{\mathbf{x}}_{n-1}$ and stops if their cosine similarity is higher than a threshold s_{prev} . We set both s_{eot} and s_{prev} to 0.9.

2.3.2 Diffusion-based LCM

Diffusion-based LCMs are generative latent variable models that learn a model distribution p_{θ} approximating a data distribution q . Similar to the BASE-LCM, we model the diffusion LCMs as auto-regressive models that generate concepts in a document, one at a time. The model distribution is thus expressed at each position n of the sequence as $p_{\theta}(\mathbf{x}_n | \mathbf{x}_{<n})$ *i.e.*, the generation of the next concept is conditioned on the preceding context.

In what follows we use a superscript for the denoising/diffusion step ($t \in [0, 1]$) and a subscript (n) for indexing the sequence of concepts. We simplify for a given n the conditional model distribution $p_{\theta}(\mathbf{x}_n^0 | \mathbf{x}_{<n}^0)$ as $p_{\theta}(\mathbf{x}^0)$, and the conditional data distribution $q(\mathbf{x}_n^0 | \mathbf{x}_{<n}^0)$ as $q(\mathbf{x}^0)$.

Diffusion models involve two processes: a *forward* noising process and a *reverse* denoising process (Ho et al., 2020; Song et al., 2020):

Forward process and noise schedule The forward process is a Gaussian diffusion process characterized by the marginal distribution $q(\mathbf{x}^t | \mathbf{x}^0)$, given for every timestep $t \in [0, 1]$ as:

$$q(\mathbf{x}^t | \mathbf{x}^0) := \mathcal{N}(\alpha_t \mathbf{x}^0, \sigma_t^2 \mathbf{I}). \quad (7)$$

With the reparameterization trick, we can sample from this marginal distribution via:

$$\mathbf{x}^t = \alpha_t \mathbf{x}^0 + \sigma_t \boldsymbol{\epsilon} \quad \text{where } \boldsymbol{\epsilon} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}) \quad (8)$$

We use a variance-preserving forward process (Karras et al., 2022) for which we have:

$$\alpha_t^2 = \text{sigmoid}(\lambda_t), \quad \sigma_t^2 = \text{sigmoid}(-\lambda_t) = 1 - \text{sigmoid}(\lambda_t), \quad \lambda_t = \log(\alpha_t^2 / \sigma_t^2), \quad (9)$$

where λ_t is the log signal-to-noise ratio (log-SNR) for timestep t .

The noise schedule is a strictly monotonically decreasing function f_{λ} that maps from the timestep $t \in [0, 1]$ to a log-SNR level: $\lambda_t = f_{\lambda}(t)$.

It is common in previous work to also define the noise schedule based on a discrete variance schedule $(\beta_0, \dots, \beta_T)$. This stems from the formulation of the forward process as a discrete-time Markov chain that gradually adds Gaussian noise to the data according to said variance schedule:

$$q(\mathbf{x}^{1 \dots T} | \mathbf{x}^0) := \prod_{t=1}^T q(\mathbf{x}^t | \mathbf{x}^{t-1}), \quad q(\mathbf{x}^t | \mathbf{x}^{t-1}) := \mathcal{N}(\mathbf{x}^t; \sqrt{1 - \beta_t} \mathbf{x}^{t-1}, \beta_t \mathbf{I}), \quad (10)$$

where to simplify the notation, \mathbf{x}^t is short for $\mathbf{x}^{t/T}$ as the timesteps are now discretized.

From the variance schedule $(\beta_t)_t$, the noise schedule can be expressed as:

$$\alpha_t^2 = \prod_{s=1}^t (1 - \beta_s). \quad (11)$$

Following [Kingma and Gao \(2024\)](#), for any given noise schedule, we visualize the distribution over noise levels $p(\lambda) = -dt/d\lambda$ in order to characterize how much time we are spending at every noise level during training.

In this work, we consider three types of noise schedules:

Cosine. The schedule formulated in [Nichol and Dhariwal \(2021\)](#) as:

$$\alpha_t^2 = f(t)/f(0), \text{ where } f(t) = \cos^2\left(\frac{t+s}{1+s} \cdot \frac{\pi}{2}\right), \text{ where } s = 0.008. \quad (12)$$

Quadratic. The schedule introduced in [Ho et al. \(2020\)](#) where the variances $(\beta_t)_t$ are set to constants increasing quadratically from β_0 to β_1 .

$$\beta_{t/T} = \left(\sqrt{\beta_0} + \frac{t}{T} \cdot (\sqrt{\beta_1} - \sqrt{\beta_0})\right)^2. \quad (13)$$

Sigmoid. We introduce in this work, the *sigmoid* schedule as a means to study the impact of the SNR distribution on the training of our models. The schedule is parametrized by two hyper-parameters (γ, δ) and is defined as:

$$\alpha_t^2 = f(t)/f(0), \text{ where } f(t) = \text{sigmoid}(\delta - \gamma \text{logit}(t)), \quad (14)$$

where ‘‘sigmoid’’ is the sigmoid function $\text{sigmoid } x \mapsto e^x/(e^x + 1)$ and ‘‘logit’’ its inverse function $\text{logit} : x \mapsto \log(x/(1-x))$. The hyper-parameter γ controls the scale of the log-SNR distribution $p(\lambda)$ and δ its center (see [Figure 5](#)).

In all our experiments, we follow [Lin et al. \(2024\)](#) and rescale the variance schedule $(\beta_1, \dots, \beta_T)$ to enforce zero terminal SNR *i.e.*, $\beta_T = 1$.

Reverse process and objective function The joint distribution of the diffusion model $p_\theta(\mathbf{x}^{0:\dots 1})$ is called the reverse process and is defined as a Markov chain with learned Gaussian transitions starting at $p(\mathbf{x}^1) = \mathcal{N}(\mathbf{0}, \mathbf{I})$. In its discretized form:

$$p_\theta(\mathbf{x}^{0:T}) := p(\mathbf{x}^T) \prod_{t=1}^T p_\theta(\mathbf{x}^{t-1} | \mathbf{x}^t), \quad p_\theta(\mathbf{x}^{t-1} | \mathbf{x}^t) := \mathcal{N}(\mathbf{x}^{t-1}; \boldsymbol{\mu}_\theta(\mathbf{x}^t, t), \boldsymbol{\Sigma}_\theta(\mathbf{x}^t, t)), \quad (15)$$

where $\boldsymbol{\mu}_\theta$ and $\boldsymbol{\Sigma}$ are predicted statistics. $\boldsymbol{\Sigma}$ is set to to the constant $\sigma_t^2 \mathbf{I}$ (matching the transitions of the forward process). $\boldsymbol{\mu}_\theta$ can be decomposed into a linear combination of \mathbf{x}^{t-1} and a noise approximation model $\boldsymbol{\epsilon}_\theta$. This prediction method is dubbed $\boldsymbol{\epsilon}$ -prediction ([Ho et al., 2020](#); [Nichol and Dhariwal, 2021](#); [Nichol et al., 2022](#)). In this work we adopt \mathbf{x}^0 -prediction *i.e.*, we predict the noiseless state and optimize the simple reconstruction loss:

$$\mathcal{L}(\boldsymbol{\theta}) := \mathbb{E}_{t \sim \mathcal{U}(0,1)} [\omega(t) \mathcal{L}(t, \boldsymbol{\theta})], \quad \mathcal{L}(t, \boldsymbol{\theta}) := \mathbb{E}_{\mathbf{x}^0, \boldsymbol{\epsilon}} \left[\left\| \mathbf{x}^0 - \boldsymbol{\mu}_\theta(\alpha_t \mathbf{x}^0 + \sigma_t \boldsymbol{\epsilon}, t) \right\|_2^2 \right]. \quad (16)$$

Different weighting strategies for the reconstruction loss were proposed in the literature ([Ho et al., 2020](#); [Salimans and Ho, 2022](#); [Hang et al., 2023](#)). In this work, we default to the simple reconstruction loss $(\omega(t) = 1, \forall t)$ and we experiment with a clamped-SNR weighting strategy:

$$\omega(t) = \max(\min(\exp(\lambda_t), \lambda_{\max}), \lambda_{\min}), \lambda_t = \log(\alpha_t^2 / \sigma_t^2), \quad (17)$$

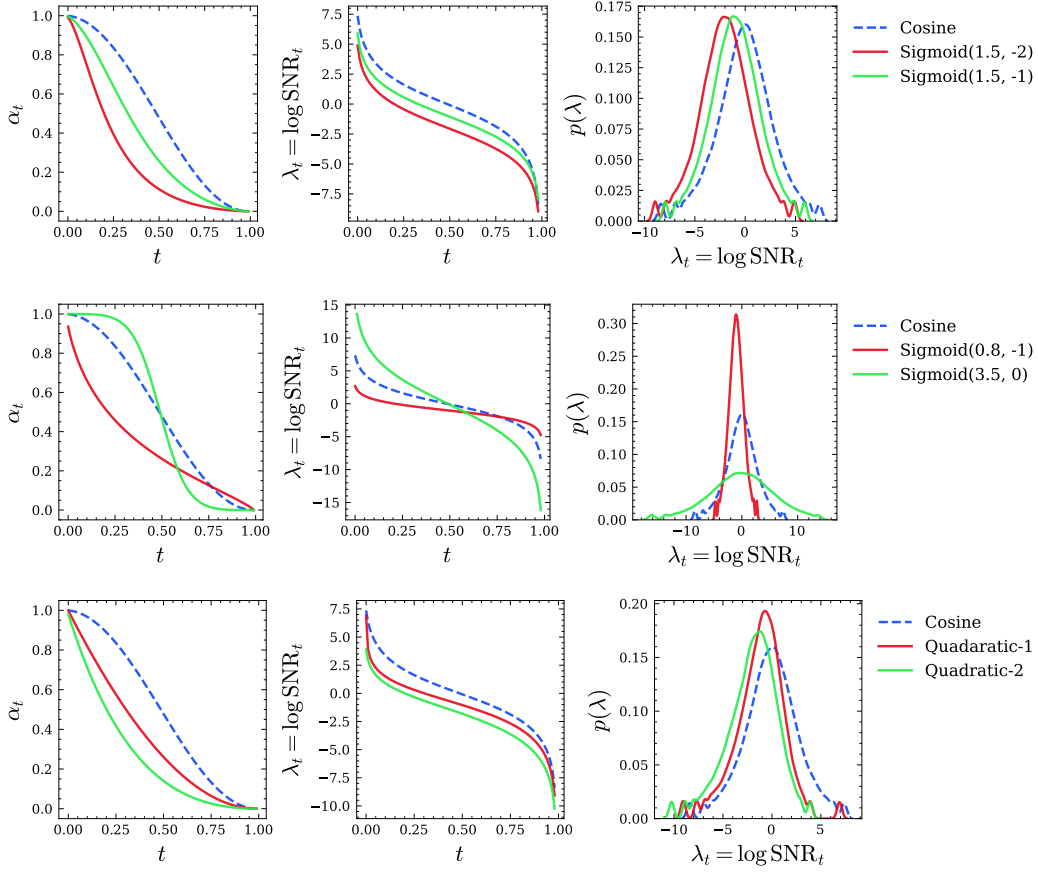


Figure 5 - Noise schedules. Illustrations of the different noise schedules explored in this work. Our default schedule being cosine. Quadratic-1 is characterized by $(\beta_0 = 0.001, \beta_T = 0.0012)$ and Quadratic-2 by $(\beta_0 = 0.02, \beta_T = 0.022)$ For each schedule we visualize the curve of $(\alpha_t)_t$ (see Equation (8)), the curve of the log-SNR and the associated distribution over noises levels $p(\lambda)$ (Kingma and Gao, 2024).

which is a generalization of Salimans and Ho (2022)’s truncated-SNR weighting and Hang et al. (2023)’s min-SNR strategy where the SNR is clamped between a min- and max-value λ_{\min} and λ_{\max} .

Additionally, we consider a weighting strategy that factors in the quality of the sample \mathbf{x}^0 . We use as sample weight a scalar $\omega(\mathbf{x}^0) \in [0, 1]$ correlated with the sample’s fragility score *i.e.*, how easy it is to reconstruct a noised sample (see Section 2.5.2). Fragile samples will be assigned a smaller weight and thus contribute less to the objective function.

$$\mathcal{L}_{\text{fragility}}(\theta) := \mathbb{E}_{t \sim \mathcal{U}(0,1), \mathbf{x}^0, \epsilon} \left[\omega(\mathbf{x}^0) \left\| \mathbf{x}^0 - \mu_{\theta}(\alpha_t \mathbf{x}^0 + \sigma_t \epsilon, t) \right\|_2^2 \right], \quad (18)$$

$$\omega(\mathbf{x}^0) = \text{sigmoid}(a \text{ fragility}(\mathbf{x}^0) + b), \quad (19)$$

where $a < 0$ and b are hyper-parameters to tune.

Classifier-free diffusion guidance for the LCM Classifier-free diffusion guidance (Ho and Salimans, 2022) consists of jointly training a conditional and an unconditional diffusion model. The resulting conditional and unconditional score estimates are combined at inference time to achieve a trade-off between sample quality and diversity. This combined score is defined as follows:

$$\nabla_x \log_{\gamma} p(x|y) = (1 - \gamma) \nabla_x \log p(x) + \gamma \nabla_x \log p(x|y), \quad (20)$$

where \mathbf{y} is the conditioning variable, in our case the sequence of preceding embeddings $(\mathbf{x}_1, \dots, \mathbf{x}_{n-1})$ when denoising \mathbf{x}_n .

The hyper-parameter γ controls the contribution of the conditional score; For $\gamma = 0$, this is equivalent to an unconditional model, and for $\gamma = 1$, it is a fully conditional model. In practice for vision models, γ is set to a value greater than 1, thus amplifying the signal from the conditioning model.

Inference At inference time, the reverse process is applied. \mathbf{x}^T is obtained by sampling a random noise from $p(\mathbf{x}^T) = \mathcal{N}(\mathbf{0}, \mathbf{I})$, and is then iteratively denoised by taking steps in the direction of the score function (*i.e.*, the direction in which the log-likelihood increases fastest). Additional noise is added during the process in order to avoid falling down into modes of the distribution.

Practically, we start from $\mathbf{x}^T \sim \mathcal{N}(\mathbf{0}, \sigma_{\text{init}}^2 \mathbf{I})$ and find that the quality of the sampled output is sensitive to the initial noise scale σ_{init} .

Although we train the model on a large number of discretized timesteps, *e.g.*, $T=100$, we only generate with a smaller number of steps, *e.g.*, $S=40$, at inference via accelerated generation processes (Song et al., 2020). We select the sample steps following the trailing method of Lu et al. (2022) as it is found to be more efficient for smaller steps S (Lin et al., 2024). That is we generate along the sampled steps $(\tau_1, \dots, \tau_S) = \text{round}(\text{flip}(\text{arange}(T, 0, -T/S)))$. During inference, we employ the classifier-free guidance rescaling technique of Lin et al. (2024) proven to alleviate the image over-exposure problem encountered in image synthesis diffusion models as the terminal SNR approaches zero. We denote with g_{scale} and g_{rescale} the guidance scale and guidance rescale factors used at inference.

Following Ning et al. (2023), we perform Epsilon-scaling at inference time as it is shown to alleviate the exposure bias problem in diffusion models. In its simplified version, this is a training-free method that consists of scaling down the over-predicted magnitude of error by a scalar λ_{eps} .

We describe in Section 2.3.3 and Section 2.3.4 two variants of diffusion LCM: ONE-TOWER and TWO-TOWER.

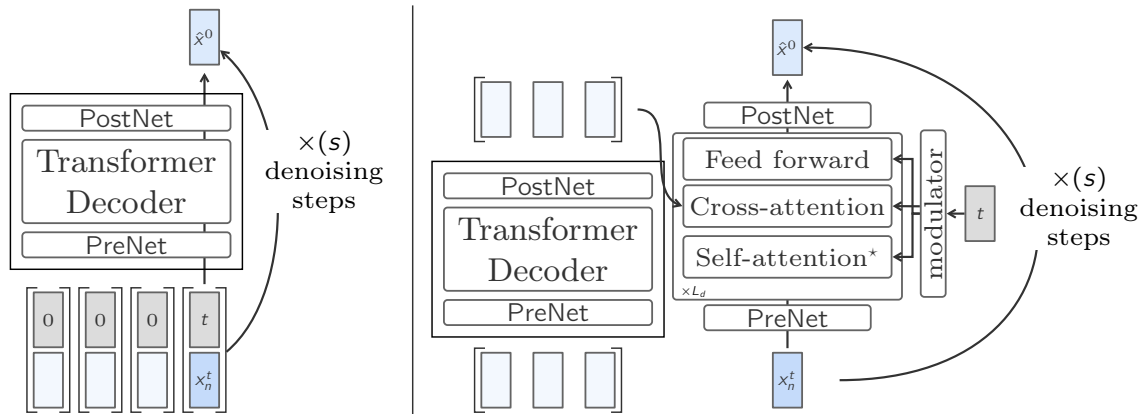


Figure 6 - Inference with diffusion-based LCMs. In the left-hand side, an illustration of the ONE-TOWER LCM and on the right-hand side an illustration of the TWO-TOWER LCM.

2.3.3 One-Tower Diffusion LCM

This model, depicted in the left panel of Figure 6, consists of a single transformer backbone whose task is to predict the clean next sentence embedding \mathbf{x}_n^0 given a noisy input \mathbf{x}_n^t , conditioned on previous clean sentence embeddings $\mathbf{x}_{<n}^0$. During training, self-attention can be dropped with a

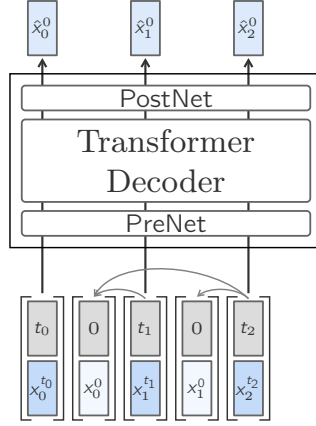


Figure 7 - Training of One-Tower diffusion LCM. Interleaving the clean and noisy embeddings and sampling different diffusion timesteps allows for efficient training.

certain probability for unconditional training. This enables classifier-free guidance at inference time (see Section 2.3.2 for details).

Each input embedding is concatenated with the corresponding diffusion timestep embedding. The learned position embeddings are added to the input vectors prior to being fed to LCM. The backbone utilizes a causal multi-head self-attention.

For efficient training, the model is trained to predict each and every sentence in a document at once. As depicted in Figure 7, during the diffusion process, the model attends to the clean sentences in the context using causal multi-head attention layers. The input is specially prepared by interleaving the noisy (blue) and clean (light blue) sentence embeddings, and the attention mask is prepared accordingly to only attend to the clean sentence embeddings (gray arrows).

2.3.4 Two-Tower Diffusion LCM

This model, depicted in the right panel of Figure 6, separates the encoding of the preceding context from the diffusion of the next embedding. A first model, labeled *contextualizer*, takes as input the context vectors $x_{<n}$ and encodes them causally *i.e.*, we apply a decoder-only Transformer with causal self-attention. The outputs of the contextualizer are then fed to a second model dubbed *denoiser*, which predicts the clean next sentence embedding \mathbf{x}_n^0 by iteratively denoising the latent $\mathbf{x}_n^1 \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$. The denoiser consists of a stack of Transformer blocks with cross-attention block to attend over the encoded context. Both the denoiser and the contextualizer share the same Transformer hidden dimension d_{model} . Each block of each Transformer layer in the denoiser (including the cross-attention layer) is modulated with adaptive layer norm (AdaLN, Perez et al. (2018); Peebles and Xie (2023)). The AdaLN modulator of TWO-TOWER regresses channel-wise scale (γ), shift (β) and residual gates (α) from the embedding of the current diffusion timestep t .

$$[\beta, \gamma, \alpha] = \text{SiLU}(\text{embed}(t))\mathbf{W}^t + \mathbf{b}, \quad (21)$$

$$\mathbf{y} = \mathbf{x} + \alpha \text{Block}((1 + \gamma) \mathbf{x} + \beta), \quad (22)$$

Following Peebles and Xie (2023) and Goyal (2017) we initialize each residual block in a Transformer layer (“Block”) with the identity function via initializing \mathbf{W} and \mathbf{b} in Equation (21) to zero. The

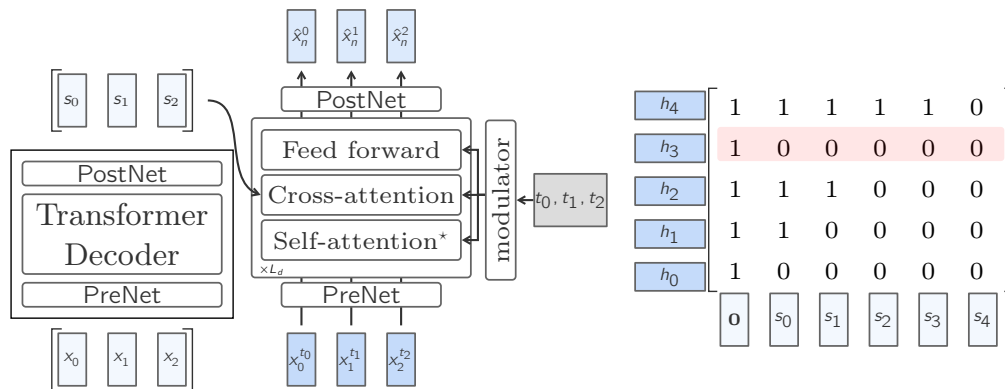


Figure 8 - Training Two-Tower diffusion LCM. On the left panel, a TWO-TOWER forward pass in training time in order to denoise multiple embeddings in parallel. On the right side panel a visualization of the denoiser’s cross-attention masks with the red highlighted row signaling a sample dropped to train the denoiser unconditionally. (h_1, \dots, h_4) denotes the sequence of intermediate representations in the denoiser right before the cross-attention layer.

diffusion timestep t is embedded using a 256-dimensional frequency embedding (Dhariwal and Nichol, 2021; Peebles and Xie, 2023) followed by a two-layer MLP with SiLU as activation function. “embed” maps to the denoiser’s hidden dimension d_{model} . The self-attention layers in the denoiser do only attend to the current position *i.e.*, we do not attend to the preceding noised context. The self-attention layers were kept for consistency with a standard Transformer block and for the possible extension of denoising multiple vectors at once.

Two-TOWER training. At training time, TWO-TOWER’s parameters are optimized for the next-sentence prediction task on unsupervised sequences of embeddings. The causal embeddings from the contextualizer are shifted by one position in the denoiser and a causal mask is used in its cross-attention layers. A zero vector is prepended to the context vectors to enable the prediction of the first position in the sequence (see Figure 8). To train the model both conditionally and unconditionally in preparation for inference with classifier-free guidance scaling, we drop random rows from the cross-attention mask with a rate of p_{cfg} and denoise the corresponding positions with only the zero vector as context.

2.3.5 Quantized LCM

Two major approaches currently stand to deal with continuous data generation in the image or speech generation fields: one is diffusion modeling, the other is learning quantization of the data before modeling on top of these discrete units.

In addition, the text modality remains discrete, and despite dealing with continuous representations in the SONAR space, all possible text sentences (of less than a given number of characters) are a cloud of points rather than a real continuous distribution in the SONAR space. These considerations motivate the exploration of quantization of SONAR representations and then modeling on these discrete units to address the next sentence prediction task. Finally, following such an approach enables the natural use of temperature, top-p or top-k sampling, to control the level of randomness and diversity in the sampling of the next sentence representation.

In this section, we learn residual quantizers for the SONAR space, and then build a Quantized LARGE CONCEPT MODEL based on these discrete units. We tried to come up with an architecture as close as the diffusion LCM models, to be able to compare approaches.

Quantization of SONAR space. We use Residual Vector Quantization (RVQ; Zeghidour et al. (2021)) as a coarse-to-fine quantization technique to discretize SONAR representations. Vector quantization maps continuous input embeddings to the nearest entry in a learnt codebook. RVQ iteratively quantize residual errors from previous quantizations using additional codebook for each iteration. We use FAISS implementation (Douze et al., 2024) which performs iterative k-means clustering of residuals. We use the Improved Residual Vector Quantization (IRVQ) method from Liu et al. (2015), with a beam size of 1 for memory efficiency. We trained the RVQ codebooks on 15 million English sentences extracted from Common Crawl using $n_{\text{codebooks}} = 64$ number of quantizers with $n_{\text{units-per-codebook}} = 8192$ units per codebook.

One property of RVQ is that the cumulative sum of centroid embeddings of the first codebooks are an intermediate coarse approximation of input SONAR vectors. In that way, we can report the evolution of auto-encoding BLEU scores with the increasing number of codebooks used to quantize SONAR embeddings, before using the SONAR text decoder to decode quantized embeddings. We notice in Figure 9 that auto-encoding BLEU consistently improves as the number of codebooks increases, reaching around 70% of the auto-encoding BLEU score achieved with continuous SONAR embeddings, when using all 64 codebooks.

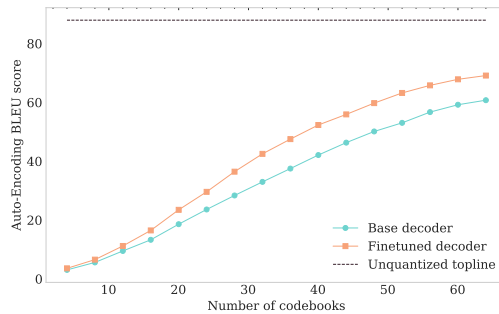


Figure 9 - Auto-encoding BLEU scores on FLORES devtest set, encoding sentences with SONAR encoder, quantizing with a varying number of codebooks, dequantizing and decoding with SONAR decoder.

Finetuning the SONAR decoder on quantized representations. We fine-tuned SONAR decoder on quantized representations to adjust it for the space created by the quantizers on 1.2M English sentences. To make the decoder more robust against residual representations from intermediate codebooks, we randomly select a codebook number $k \in \left[\frac{2}{3} \cdot n_{\text{codebooks}}, n_{\text{codebooks}}\right]$ during fine-tuning, with probability $p = 0.3$, and use the quantized representation with codebooks up to k . Figure 9 shows the improvement in auto-encoding performance when the decoder is adapted to quantized representations.

QUANT-LCM architecture. In the same spirit of diffusion LCM, we aim at coarse-to-fine generation of SONAR embeddings conditioned on left-context sentences. However, we do not follow a denoising task as in diffusion modeling, but an iterative generation of SONAR embeddings based on *intermediate quantized representations* instead. In order to generate a SONAR embedding conditioned on left-context sentences, the QUANT-LCM model starts with the *intermediate representation* as a vector filled with zeros. We iteratively add to this *intermediate representation* the predicted residual centroid embeddings. In that way, the predicted SONAR embeddings are iteratively refined based on the growing cumulative sum of centroid embeddings of first codebooks, until all codebooks

have been seen. We used the ONE-TOWER architecture for QUANT-LCM experiments even though it could be trained with TWO-TOWER architecture too. Compared to the diffusion LCM, noisy input representations are replaced with *intermediate quantized representations* and diffusion timestep embeddings as input are replaced by codebook index embeddings.

Discrete targets. Following previous work on modeling discrete units from residual quantizers (Wang et al., 2023; Rubenstein et al., 2023; Lee et al., 2022), a Quant-LCM can be trained to predict the unit from the next codebook, parameterized with a softmax output layer. For parameter efficiency, we do not use $n_{\text{codebooks}} \cdot n_{\text{units-per-codebook}}$ unique indices as discrete targets which would imply $n_{\text{codebooks}} \cdot n_{\text{units-per-codebook}}$ output dimensions, but only $n_{\text{units-per-codebook}}$ output dimensions while inputting the information of the codebook index to the model. At training time, similarly to diffusion LCM training, we randomly sample codebook index k between 1 and $n_{\text{codebooks}}$, and compute the cumulative sum of centroid embeddings of the first $k-1$ codebooks as input. We use the unit from codebook k of the target embedding as target index for cross entropy loss computation. At inference time, we iteratively predict the unit from the next codebook, get the corresponding centroid embedding and add it to the current *intermediate representation* as additional predicted residual embedding. Finally, we also enable classifier-free guidance on logits at inference time (Gafni et al., 2022) by randomly dropping left-context conditioning during training as previously described in Section 2.3.3. This modeling approach with discrete targets is dubbed QUANT-LCM-D in the following sections. The improved SONAR decoder for quantized representations is used to bridge the compression gap coming from SONAR quantization in following ablation studies when using QUANT-LCM-D.

Continuous targets. We also explored a modeling approach that predicts continuous target SONAR vectors based on left-context sentences and intermediate quantized representation of the target vector, minimizing the Mean Squared Error between prediction and target embeddings. At inference time, we can either iteratively add the closest centroid embedding based on the predicted residual $\hat{\mathbf{r}}$ or sample a centroid \mathbf{c}_i from the following distribution:

$$p(\mathbf{c}_i|\hat{\mathbf{r}}) = \frac{e^{-\beta \cdot \|\mathbf{c}_i - \hat{\mathbf{r}}\|_2}}{\sum_k e^{-\beta \cdot \|\mathbf{c}_k - \hat{\mathbf{r}}\|_2}}, \quad (23)$$

where β is a temperature hyper-parameter. This modeling approach with continuous targets is denoted with QUANT-LCM-C in the following sections.

2.4 Ablations

In this section, we delineate the ablations experiments conducted to evaluate the aforementioned LCM designs. We compare all the variants of LCMs introduced above, namely, BASE-LCM, ONE-TOWER, TWO-TOWER and QUANT-LCM.

2.4.1 Experimental setup

For our ablation study and for the sake of reproducibility, we pre-train our models on the FINEWEB-EDU dataset (Lozhkov et al., 2024). All models are configured to have approximately 1.6B trainable parameters and are pre-trained on Meta’s Research Super Cluster (RSC, Lee and Sengupta (2022)) for 250k optimization steps spanning 32 A100 GPUs with a total batch size of 229k concepts.

Models architectures. The BASE-LCM has 32 layers and a model dimension $d_{\text{model}} = 2048$ with 16 attention heads. It uses rotary position embeddings (RoPE, [Su et al. \(2024\)](#)), applies pre-normalization using RMSNorm ([Zhang and Sennrich, 2019](#)), uses the SwiGLU activation function ([Shazeer, 2020](#)) and is trained with a dropout rate of $p=0.1$.

The ONE-TOWER diffusion LCM is made of 32 transformer blocks, each made of a self-attention layer with 32 attention heads and followed by a feed-forward neural network with inner size 8192. It has a dimension d_{model} of 2048 and uses learned position embeddings. The noise scheduler is set with $T=100$ diffusion timesteps. During training, self-attention is dropped with a probability of 0.15 for unconditional training, enabling classifier-free guidance at inference time.

The TWO-TOWER diffusion LCM has 5 layers in its contextualizer and 13 layers in its denoiser. Similar to the BASE-LCM, it has 16 attention heads, a model dimension $d_{\text{model}} = 2048$, and uses SwiGLU activations and RMSNorm in both contextualizer and denoiser. The contextualizer uses RoPE for embedding positions whereas the denoiser is without positional embeddings. We use by default the cosine noise schedule with $T=100$ and train with a dropout rate of $p=0.1$. For training the model unconditionally we use a cross-attention mask dropout of rate 0.15 (see [Section 2.3.4](#)). The pre-training documents are wrapped at 128 sentences. Unless otherwise mentioned we decode with $S=40$ sample steps with a guidance scale $g_{\text{scale}} = 3$, a guidance rescaling factor of $g_{\text{rescale}} = 0.7$, an initial noise scale $\sigma_{\text{init}} = 0.6$ and epsilon-scaling with $\lambda_{\text{eps}} = 1.00045$.

The QUANT-LCM follows exactly the same architecture as the ONE-TOWER diffusion LCM, except for QUANT-LCM-D which differs only by its output dimension which is set to $n_{\text{units-per-codebook}} = 8192$ for softmax computation. For single sentence prediction tasks, we use $\text{top}_k = 1$ and $g_{\text{scale}} = 2$ for QUANT-LCM-D and $\text{top}_k = 1$ with $g_{\text{scale}} = 3$ for QUANT-LCM-C, while for multi-sentence generation tasks we used temperature of 1, $\text{top}_k = 3$, $g_{\text{scale}} = 1$, for QUANT-LCM-D and temperature of 0.005, $\text{top}_k = 5$, $g_{\text{scale}} = 1.5$ for QUANT-LCM-C, as higher guidance or lower temperature setups led to repeated generated sentences.

Dataset	#Docs	#LLAMA2 Tokens			#Sentences			Total sentences
		Q1	Q2	Q3	Q1	Q2	Q3	
ROC-STORIES (dev)	2000	48	57	64	5	5	5	10K
ROC-STORIES (test)	1871	50	56	62	5	5	5	9.4K
C4 (dev)	1000	136	288	577	6	12	24	20.6K
C4 (test)	1000	133	282	599	6	11	25	21.9K
WIKIPEDIA-EN (dev)	1000	146	332	736	5	10	23	21.1K
WIKIPEDIA-EN (test)	1000	147	312	673	5	9	21	19.1K
GUTENBERG (dev)	55	10297	15934	22259	328	530	687	216.2K
GUTENBERG (test)	61	10752	15204	23414	325	457	735	562.2K

Table 2 - Statistics of the pre-training evaluation datasets. For each subset we report the number of documents, the total number of sentences and document lengths quartiles in sentences and in LLAMA2 tokens for reference.

Pre-training evaluation. Pre-trained token-level language models are typically evaluated with perplexity: a measure of how well each next token is predicted given a teacher-forced (*i.e.*, ground truth) prefix of the document. In a similar spirit, we evaluate pre-trained LCMs in a teacher-forced mode. But as they cannot produce the probability explicitly, we resort to a custom set of metrics of the quality of next sentence prediction.

Each pre-trained model is initially evaluated on the quality of its predicted next sentence $\hat{\mathbf{x}}_n$ given a ground truth context $\mathbf{x}_{<n}$. Practically, for a given document $\mathbf{x}_{1:N}$, we run the LCM inference in teacher-forcing mode and evaluate the following metrics:

- **L2 distance** (ℓ_2). Euclidean distance in the SONAR space between the predicted embedding $\hat{\mathbf{x}}_n$ and the ground truth continuation \mathbf{x}_n : $\ell_2 := \|\hat{\mathbf{x}}_n - \mathbf{x}_n\|^2$.
- **Round-trip L2 distance** (ℓ_{2-r}). Euclidean distance in the SONAR space between the re-encoded sentence generated from the predicted embedding and the ground truth continuation \mathbf{x}_n ,

$$\ell_{2-r} := \|\text{encode}(\text{decode}(\hat{\mathbf{x}}_n)) - \mathbf{x}_n\|^2.$$

Since an LCM can predict an embedding outside of the distribution of real embeddings (obtained by encoding natural sentences), the SONAR decoder might shift these embeddings to the nearest plausible embeddings subspace. The ℓ_{2-r} metric is introduced to capture the shift in embeddings after decoding them into text then re-embedding them again in the SONAR space. The more the generated embeddings are out-of-distribution, the higher the delta between ℓ_{2-r} and ℓ_2 would be.

- **Contrastive accuracy** (CA). The ratio of embeddings in a batch that are further away (in terms of ℓ_2) from the predicted embedding $\hat{\mathbf{x}}_n$ than the ground truth \mathbf{x}_n (for each n , we exclude \mathbf{x}_n and its two neighboring ground truth embeddings from the comparison). This metric naturally assigns higher penalty for large ℓ_2 values in the regions with high density of sentence embeddings.
- **Paraphrasing** (PAR). The maximum cosine similarity (CS) between the generated embedding $\hat{\mathbf{x}}_n$ and the context embeddings $\mathbf{x}_{<n}$, normalized by the score of the ground truth sentence. Thus, $\text{PAR} = \max_{m < n} \text{CS}(\hat{\mathbf{x}}_n, \mathbf{x}_m) / \max_{m < n} \text{CS}(\mathbf{x}_n, \mathbf{x}_m)$. The goal of this metric is to capture if the model is simply copying or paraphrasing a sentence from the context more (> 1) or less (< 1) than the reference sentence.
- **Mutual information** (MI). This metric of text coherence evaluates the mutual information between the next predicted sentence $\hat{s}_n = \text{decode}(\hat{\mathbf{x}}_n)$ and the previous $k = 10$ ground truth sentences by computing the difference between the unconditional perplexity of \hat{s}_n and its perplexity conditioned on the prompt:

$$\text{MI} = \frac{1}{|\hat{s}_n|} (\log \text{p}_{\text{LM}}(\hat{s}_n) - \log \text{p}_{\text{LM}}(\hat{s}_n | s_{n-k:n-1})).$$

We estimate the perplexity with a small language model, GPT-2 (Radford et al., 2019). We prepend a newline symbol to \hat{s}_n , so that a probability could be assigned to its first token, and we compute the average mutual information per-token by normalizing it with $|\hat{s}_n|$, the length of \hat{s}_n in tokens. When averaging MI over a dataset, $|\hat{s}_n|$ are used as weights.

Pre-training evaluation data. The pre-training evaluation is performed on sampled subsets from four corpora covering different domains: ROC-STORIES (Mostafazadeh et al., 2016), C4 (Raffel et al., 2019), WIKIPEDIA-EN (English Wikipedia dump) and GUTENBERG. We sample two distinct subsets (dev and test) from each corpus, we use the dev split for tuning inference hyper-parameters and report the results on the test splits. The statistics of the evaluation corpora are presented in Table 2.

The results of the pre-training evaluation are presented in Table 3.

First, diffusion-based LCM and QUANT-LCM variants have similar ℓ_2 and ℓ_{2-r} scores despite an important difference in their learning objectives. The only model that shows substantially lower ℓ_2 score is the BASE-LCM. This is expected since BASE-LCM effectively optimizes ℓ_2 score during training. Yet, ℓ_{2-r} score is not improved compared to other models. This could be explained by the

Model	ROC-STORIES					C4				
	ℓ_2	ℓ_{2-r}	PAR	CA	MI	ℓ_2	ℓ_{2-r}	PAR	CA	MI
BASE-LCM	0.177	0.237	1.847	72.4%	0.062	0.204	0.261	1.964	69.1%	-0.105
ONE-TOWER	0.236	0.236	1.939	80.2%	0.977	0.279	0.273	2.239	77.1%	1.110
TWO-TOWER	0.233	0.231	2.088	80.6%	1.137	0.265	0.261	2.265	75.4%	1.134
QUANT-LCM-C	0.236	0.237	1.683	76.0%	0.610	0.279	0.283	2.013	77.2%	0.715
QUANT-LCM-D	0.240	0.246	1.871	81.1%	0.682	0.270	0.282	1.808	75.0%	0.359

Model	WIKIPEDIA-EN					GUTENBERG				
	ℓ_2	ℓ_{2-r}	PAR	CA	MI	ℓ_2	ℓ_{2-r}	PAR	CA	MI
BASE-LCM	0.229	0.283	1.770	69.6%	0.071	0.207	0.264	1.780	67.8%	-0.184
ONE-TOWER	0.324	0.311	2.087	80.9%	1.202	0.284	0.281	2.051	75.1%	0.725
TWO-TOWER	0.307	0.297	2.079	78.8%	1.307	0.267	0.267	2.077	73.0%	0.684
QUANT-LCM-C	0.306	0.317	1.842	79.5%	0.744	0.269	0.281	1.774	72.1%	0.419
QUANT-LCM-D	0.295	0.311	1.592	76.0%	0.323	0.276	0.290	1.599	72.0%	0.153

Table 3 - Comparing architectures. Pre-training evaluation results on the four select corpora. For each subset, we report ℓ_2 (L2 distance in SONAR space), ℓ_{2-r} (round-trip L2 distance after decoding and re-encoding the generated embeddings), PAR (similarity to preceding embeddings) and CA (contrastive accuracy)

fact that when many plausible next sentence continuations are possible, BASE-LCM generates their average in SONAR space (instead of sampling one of plausible modes) which may not correspond to any relevant point in SONAR space. This hypothesis is also highlighted by the poor BASE-LCM performance in term of CA and MI scores.

We do not notice any significant difference in CA scores between diffusion LCMS and QUANT-LCM variants. MI scores, on the contrary, are consistently higher for diffusion-based models compared to QUANT-LCM. At the same time, diffusion LCMS tend to paraphrase more the context in the generated embeddings, which also correlates with an increased MI score. Still, QUANT-LCM variants significantly outperform BASE-LCM on MI metric. Now comparing the different variants, QUANT-LCM-C outperforms QUANT-LCM-D modeling variant: one hypothesis is that predicting codebook indices with cross-entropy loss is harder than MSE objective where QUANT-LCM-C can more easily learn combination of left-context vectors for next sentence embedding.

For diffusion LCMS, we don’t observe any consistent difference between ONE-TOWER and TWO-TOWER when looking across all metrics and datasets. Note that overall, to tackle the next sentence prediction task in the SONAR space, diffusion-based methods give clearly better results compared to all other models.

Instruction-tuning evaluation. Subsequently, the pre-trained models are instruction-tuned on the stories subset of COSMOPEDIA (Ben Allal et al., 2024) and are evaluated on a held-out subset of COSMOPEDIA itself. We aim with this finetuning to evaluate the ability of the models to follow instructions and generate consistent stories.

For the sake of comparison, we trained a small Llama (Touvron et al., 2023) on the same training data (FINEWEB-EDU) and finetuned it on COSMOPEDIA. This model has 24 transformer layers, each with 16 attention heads and a model dimension of 2048 for a total of 1.4B parameters. This model will be referred to as SMALLAMA.

We evaluate the following metrics:

- **ROUGE-L** (R-L). ROUGE-L (F-measure) (Lin, 2004) between the generated and reference stories.
- **Coherence** (Coherence). This reference-free metric is computed with a bidirectional transformer model fine-tuned by Jwalapuram et al. (2022) to assign higher scores to positive “natural” documents than to negative examples with permuted sentences. For reporting, we normalize it with a sigmoid (with a temperature 3.0, empirically set to make the scores of “certainly incoherent” documents close to 0 and those of “certainly coherent” documents close to 1).

Model	R-L \uparrow	Coherence \uparrow
BASE-LCM	23.69	0.482
ONE-TOWER	33.40	0.968
TWO-TOWER	33.64	0.938
QUANT-LCM-C	30.87	0.847
QUANT-LCM-D	28.01	0.704
SMALLAMA	34.88	0.984

Table 4 - Comparing architectures. Instruction-tuning evaluation results. For each model we score the generated stories on the held-out test prompts and report R-L (ROUGE-L) scores.

The scores in terms of R-L and Coherence of the finetuning evaluation are presented in Table 4. Those quantitative results are in line with the pretraining evaluation ones. Both R-L and Coherence scores correlate with the model ordering based from MI scores, mainly QUANT-LCM is outperformed by diffusion-based models, and both outperform BASE-LCM by a large margin.

We also note that SMALLAMA outperforms the LCMS on this downstream task on both metrics. It is well known that LLMs produce very fluent outputs, that explains the higher ROUGE-L score. We also note that the ONE-TOWER and TWO-TOWER produce coherent outputs, on par with the SMALLAMA outputs.

2.4.2 Importance of the diffusion inference hyper-parameters

In this section we study the effect of different inference hyper-parameters on the quality of the generated text. To this end, we generate outputs for the C4 test split with the TWO-TOWER LCM model above, while varying the following hyper-parameters: the guidance scale g_{scale} , the initial noise scale σ_{init} , and the number of inference sample steps S . We score the generations following the same protocol above and report the results in Figure 10. We note that as we increase the guidance scale, the mutual information between the prefix and the generated suffix increases, and so does paraphrasing as we are paying more attention to the conditioning context variables. In the opposite direction of the mutual information, the ℓ_2 distance from the ground truth continuation increases as the model prefers to stay close to the prefix. Regarding the initial noise scale, we observe that values between 0.5 and 0.7 achieve the best MI score. In particular, the generated individual sentences are usually longer with a higher σ_{init} . The ℓ_2 distance on the other hand does not reflect this trend. Lastly, we increase the number of inference steps and measure the mutual information (MI) of the generated texts. With more inference steps, we can improve the prefix-suffix mutual information, but there is diminishing returns to increasing the inference cost with little qualitative improvement.

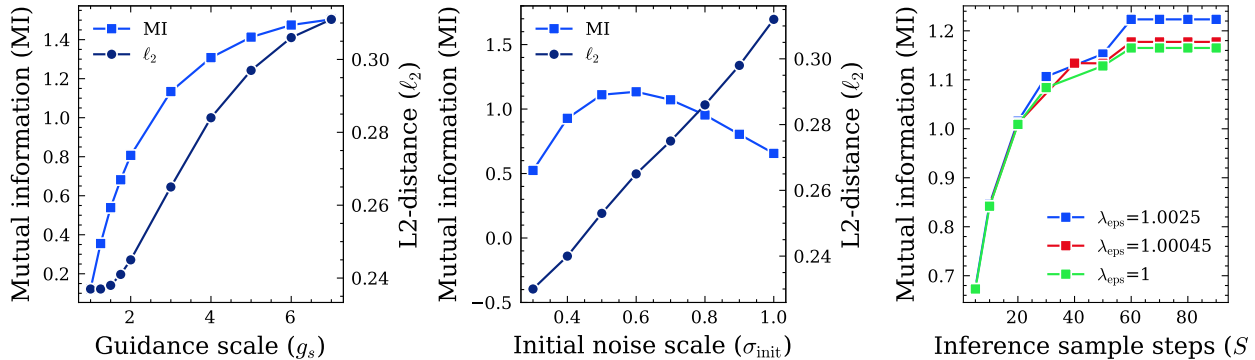


Figure 10 - Importance of inference hyper-parameters. The first panel shows the quality of the generated output measured with MI and ℓ_2 as we vary the guidance scale g_{scale} with fixed $\sigma_{init} = 0.6$ and $S = 40$. The second panel varies the initial noise scale σ_{init} with fixed guidance $g_{scale} = 3$ and $S = 40$. The third panel varies the inference steps S while holding the guidance scale $g_{scale} = 1.5$ and $\sigma_{init} = 0.6$ fixed. We consider 3 values for λ_{eps} to see the impact of epsilon-scaling in the regime of large inference steps.

2.4.3 Studying the noise schedules

In this section, we compare different TWO-TOWER diffusion LCMs trained with different noise schedules, namely:

Cosine. Our default cosine noise schedule.

Quadratic. The first quadratic schedule (Quadratic-1) has $\beta_0 = 0.001$ and $\beta_T = 0.0012$, whereas the second quadratic schedule (Quadratic-2) has $\beta_0 = 0.02$ and $\beta_T = 0.022$.

Sigmoid. Four different sigmoid schedules with with $(\alpha, \beta) \in \{(1.5, -1), (1.5, -2), (0.8, -1), (3.5, 0)\}$.

All schedules are set with the default $T=100$. For the exact description of each noise schedule refer to Section 2.3.2 and Figure 5. We selected the quadratic schedule as a commonly used schedule for reference with two configurations, Quadratic-1 closer to the cosine schedule and Quadratic-2 with more weight given to lower log-SNR. The selected sigmoid schedules with $\delta = 1.5$ are configured with $\gamma = -1$ and $\gamma = -2$, $\gamma = -2$ being slightly shifted to the left on the log-SNR distribution *i.e.*, more weight to lower log-SNR regimes. We then change the δ parameter of the sigmoid schedule to choose $(\delta = 0.8, \gamma = -1)$ for a peaked distribution of log-SNR around -1 and $(\delta = 3.5, \gamma = 0)$ for a flat distribution over noise levels.

We follow the experimental setup described in Section 2.4.1 and report the results of the pre-training evaluation in Table 5. Both quadratic schedules achieve a better MI score while the wide sigmoid schedule $(\delta, \gamma) = (3.5, 0)$ achieves the highest accuracy CA on both C4 and WIKIPEDIA-EN. To further understand the differences between those schedules we re-evaluate on C4 while varying the guidance scale g_{scale} . The results in Figure 11 confirm that the wide sigmoid schedule $(\delta = 3.5, \gamma = 0)$ has a much higher contrastive accuracy across the board (rightmost panel) while closely matching the cosine schedule in terms of mutual information (MI). This schedule being trained on a wider spectrum of log-SNR learns to contrast more than to regress. Contrast this behavior to the peaked sigmoid schedule $(\delta = 0.8, \gamma = -1)$ where the model focuses on regressing to the target (lower ℓ_2), akin to a BASE-LCM, resulting in a lower contrastive accuracy.

Model	C4					WIKIPEDIA-EN				
	ℓ_2	ℓ_{2-r}	PAR	CA	MI	ℓ_2	ℓ_{2-r}	PAR	CA	MI
Cosine	0.265	0.261	2.265	75.4%	1.134	0.307	0.297	2.079	78.8%	1.307
Quadratic-1	0.268	0.264	2.341	75.7%	1.252	0.309	0.300	2.202	79.1%	1.409
Quadratic-2	0.270	0.265	2.320	76.2%	1.252	0.312	0.303	2.185	79.7%	1.399
Sigmoid(1.5, -1)	0.257	0.259	2.226	74%	1.083	0.298	0.292	2.110	77%	1.271
Sigmoid(1.5, -2)	0.277	0.267	2.291	77.2%	1.173	0.321	0.303	2.179	80.3%	1.308
Sigmoid(0.8, -1)	0.252	0.255	2.053	70.6%	0.936	0.285	0.283	1.883	71.7%	1.127
Sigmoid(3.5, 0)	0.307	0.265	2.347	80.3%	1.154	0.347	0.303	2.187	83.7%	1.288

Table 5 - Comparing noise schedules. Results of the pre-training evaluation on two corpora, C4 and WIKIPEDIA-EN.

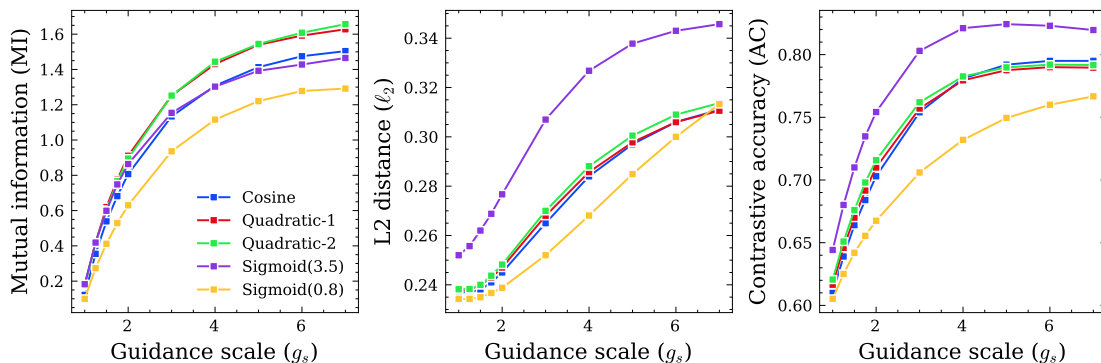


Figure 11 - Comparing noise schedules. The prefix-suffix mutual information (MI), the ℓ_2 distance and contrastive accuracy (CA) scores of evaluated C4 documents while varying the guidance scale (g_{scale}) under different schedules.

2.4.4 Studying the loss weighting strategies

In this section we compare the baseline TWO-TOWER diffusion LCM trained with the simplified objective (*i.e.*, $\omega(t) = 1, \forall t$) to models trained with the clamped-SNR weighting strategy. We consider two sets of $(\lambda_{\min}, \lambda_{\max})$: $(\lambda_{\min}, \lambda_{\max}) = (0, 10)$ and $(\lambda_{\min}, \lambda_{\max}) = (0.001, 5)$. All models in this section are trained with the cosine noise schedule.

Fragility as a sample weighing strategy As introduced in Equation (19), we train in this section a TWO-TOWER LCM model with loss terms weighted by the sample’s fragility.

$$\omega(\mathbf{x}^0) = \text{sigmoid}(a \text{ fragility}(\mathbf{x}^0) + b) \quad (24)$$

Given that estimating the fragility score of each sample as defined in Section 2.5.2 can be very costly, we resort to training a simple MLP (3-layers) on 50M sampled sentences to approximate these fragility scores. This model is referred to as F. Henceforth, the sample weight is:

$$\omega(\mathbf{x}^0) = \text{sigmoid}(a F(\mathbf{x}) + b) \quad (25)$$

The two hyper-parameters (a, b) are chosen so that extremely fragile sentences contribute less to the loss ($\omega(\mathbf{x}^0) \approx 0$), and so that sample weights should increase smoothly as sample robustness

improves. Figure 12 plots the sample weights distribution evaluated on a pre-training dataset with $a = -4$ and $b = 3.5$.

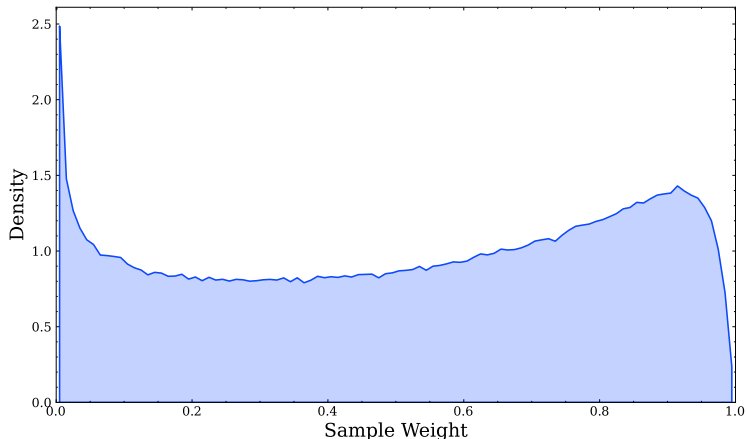


Figure 12 - Resulting distribution of fragility sample weights $\omega(\mathbf{x}^0)$ with $\omega(\mathbf{x}^0) = \text{sigmoid}(-4 F(\mathbf{x}) + 3.5)$.

We follow the experimental setup described in Section 2.4.1 and report the results of the pre-training evaluation in Table 6. We observe that weighting with clamped SNRs does not improve the quality of generated texts as measured with our pre-training evaluation metrics. When it comes to the fragility-aware weighting strategy, we observe an improvement in the contrastive accuracy of the model. In the remainder of this work we default to the simplified training objective ($\omega(t) = 1, \forall t$).

Model	C4					WIKIPEDIA-EN				
	ℓ_2	ℓ_{2-r}	PAR	CA	MI	ℓ_2	ℓ_{2-r}	PAR	CA	MI
Baseline $\omega(t) = 1$	0.265	0.261	2.265	75.4%	1.134	0.307	0.297	2.079	78.8%	1.307
SNR (0,10)	0.280	0.264	2.334	74.8%	1.107	0.320	0.296	2.102	77.9%	1.212
SNR (0.001,5)	0.266	0.261	2.269	73.4%	1.094	0.304	0.291	2.007	76.6%	1.295
Fragility	0.2815	0.273	2.306	76.5%	1.103	0.321	0.308	2.118	80.7%	1.193

Table 6 - Comparing weighting strategies. Results of the pre-training evaluation on two corpora, C4 and WIKIPEDIA-EN.

2.5 Analysis

2.5.1 Inference efficiency of LCMs

We compare in this section the inference computational cost of the TWO-TOWER LCM to that of a vanilla LLM as a function of the total length in tokens of the prompt and output combined. We chose the theoretical number of FLOPs, independent of any specific optimization. These optimizations are generic to the transformer architecture and also apply to our LCM. We include in this comparison two configurations of LCMs; the 1.6B used in the previous ablation studies and a 7B model we scale up to in the following sections. For both LCMs, we estimate the inference cost with inference sample steps $S = 40$. Given the quadratic complexity of the attention mechanism in transformers, the complexity sharply increases with the context size (see upper right corner of Figure 13’s left panel). The complexity of the LCM depends on how the context is sentencized: a context length of 200 tokens split into 10 sentences (20 tokens each) will incur a higher cost than the same 200 tokens split

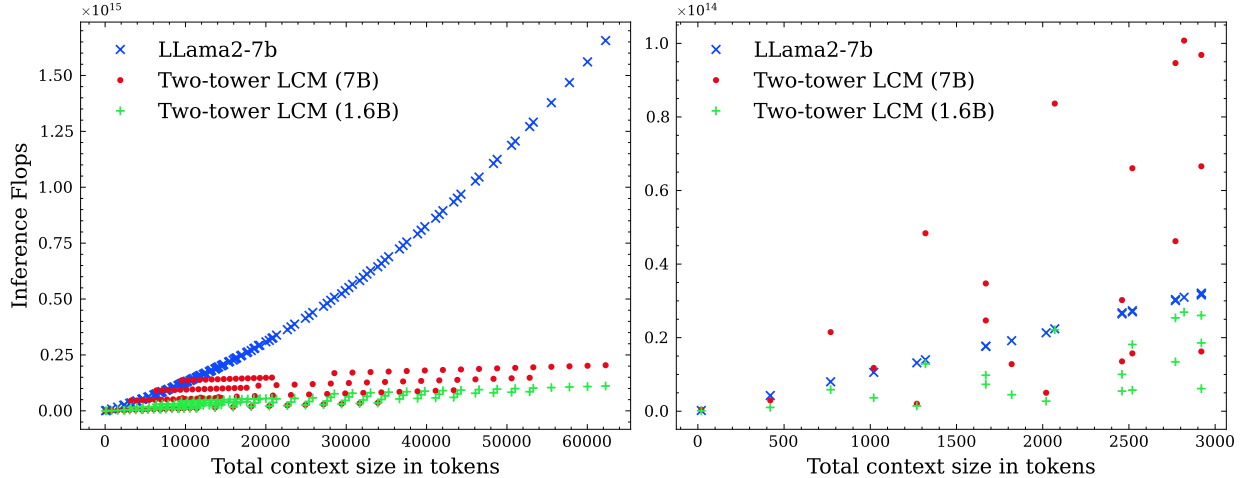


Figure 13 - Theoretical inference Flops of LCMs and LLMs. We evaluate the inference flops for different text lengths (in LLAMA2 tokens) with a variable average sentence length. Only extremely short sentences (≤ 10 tokens) favor LLMs.

into 5 sentences (40 tokens each). We account for this by computing the cost on a range of sentence lengths but report the total context size on the x-axis (context size = sentence length \times number of sentences). The LCM shows substantially better scalability with respect to increasing context size. The inference computational cost of the LCM includes the three steps of (1) encoding into SONAR, (2) LCM prediction in the sentence space then (3) decoding with a SONAR decoder. The inference cost of LCMS varies significantly depending on the average length in tokens per sentence. For extremely short sentences (less than 10 tokens), an LLM is more computationally efficient (see lower left corner of Figure 13’s right panel).

2.5.2 Fragility of SONAR space

When we perform modeling in a latent space, we primarily rely on the induced geometry ($L2$ -distance). However, the homogeneous Euclidean geometry of any latent representation will not perfectly match the underlying text semantics. This is evidenced by the fact that a small perturbation in the embedding space may result in a drastic loss of semantic information after decoding. We dub such embeddings “*fragile*”. For this reason, we aim to quantify the fragility of semantic embeddings (namely SONAR codes) to understand the quality of the LCM training data and how this fragility can hinder the LCM training dynamics.

Given a text fragment w and its SONAR code $\mathbf{x} = \text{encode}(w)$, we define the fragility of w as:

$$\text{fragility}(w) := -\mathbb{E}_{\alpha \sim \mathcal{U}([0,1]), \epsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I})} [\text{score}(w, w_{\alpha, \epsilon})], \quad (26)$$

$$\mathbf{x}_{\alpha, \epsilon} = \text{denormalize}(\sqrt{1 - \alpha} \text{normalize}(\mathbf{x}) + \sqrt{\alpha} \epsilon), \quad (27)$$

$$w_{\alpha, \epsilon} = \text{decode}(\mathbf{x}_{\alpha, \epsilon}), \quad (28)$$

where normalize and denormalize are the normalization and denormalization operators introduced in Equation (4) with the goal of making \mathbf{x} ’s coordinates scale-independent. The “encode“ operation maps text fragments into SONAR space, and the “decode“ operation produces a text fragment from a given vector in the SONAR space. For each α in $[0, 1]$, $\mathbf{x}_{\alpha, \epsilon}$ is the perturbed version of \mathbf{x} where a noise vector of variance α is linearly combined with \mathbf{x} . The perturbed vector is then decoded into a

text fragment $w_{\alpha,\epsilon}$. This perturbation is similar to the variance-preserving noising used in diffusion LCMs (see [Section 2.3.2](#)).

The “score” operator in [Equation \(26\)](#) is set to be a semantic similarity metric comparing the perturbed text $w_{\alpha,\epsilon}$ to the original w . We considered the following options:

- **Auto-Encoding BLEU.** $\text{score}(w, w_{\alpha,\epsilon}) = \text{BLEU}(w, w_{\alpha,\epsilon})$.
- **External cosine similarity.** Provided an external text encoder (read unrelated to SONAR) that encodes a text fragment w into $\text{encode}_{\text{ext}}(w)$ $\text{score}(w, w_{\alpha,\epsilon}) = \text{CS}(\text{encode}_{\text{ext}}(w), \text{encode}_{\text{ext}}(w_{\alpha,\epsilon}))$, where CS is the cosine similarity measure. Compared to Auto-Encoding BLEU, this method is typically more robust to paraphrasing.

Finetuned robust decoder. To serve as a testbed for our fragility analysis, a new SONAR decoder for English text is finetuned on a sample of our pre-training data. In order to improve the decoder’s robustness to imperfectly generated embeddings from the LCM, we follow [Equation \(27\)](#) and add random noise vectors to SONAR embeddings during training. As reported in [Table 7](#), the finetuned SONAR decoder exhibits stronger performance across a range of public corpora.

Model	FLORES	CNN DAILYMAIL	GUTENBERG	C4
BASE SONAR decoder	79.5	75.9	70.5	75.7
FINETUNED SONAR decoder	88.0	87.6	85.6	87.5

Table 7 - Comparing SONAR decoders. Raw reconstruction performance of our base SONAR decoder vs. the new decoder trained with noised embeddings. Scores are Auto-Encoding BLEU on random subsets of 10k sentences from each dataset, except for FLORES where we use the entire dev split.

Fragility study. We sample 50M random text fragments, and for each sample we generate 9 perturbations corresponding to different noise levels $\alpha \in [0.1, 0.2, \dots, 0.9]$. For the external cosine similarity metric we use MGTE as external encoder ([Zhang et al., 2024](#)).

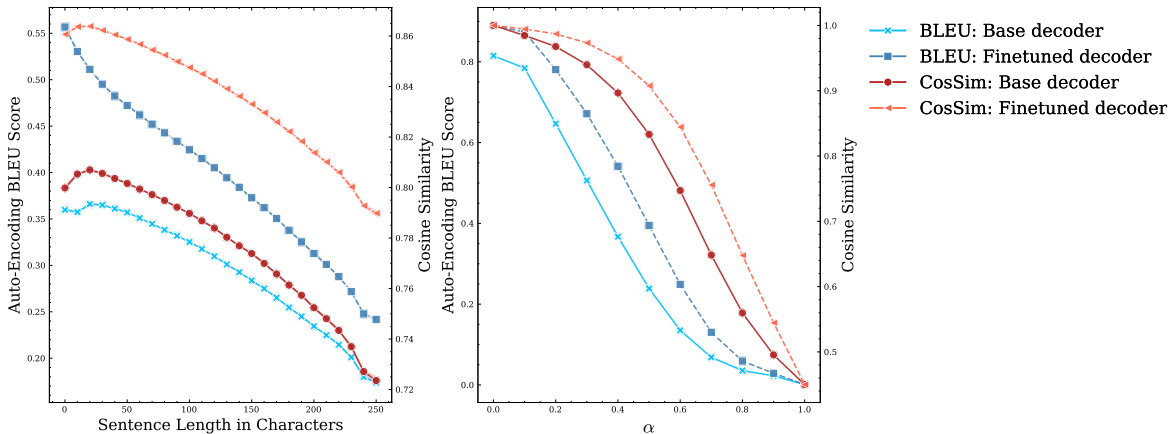


Figure 14 - Fragility scores. Auto-Encoding BLEU and external cosine similarity. In the left-hand panel as a function of the text length (α -averaged) and in the right-hand panel as a function of the noise variance α .

We depict in the right panel of [Figure 14](#) the curves of both score functions with respect to the noise level α . We observe that BLEU scores decrease faster than the cosine similarity. Most importantly, fragility scores are sensitive to the choice of the decoder. In particular, both Auto-Encoding BLEU

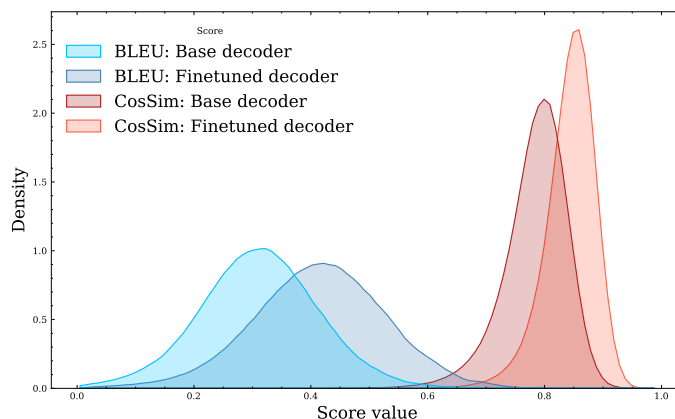


Figure 15 - Auto-Encoding BLEU scores and cosine similarity (α -averaged) distributions.

and cosine similarity scores decrease at a markedly slower rate for the FINETUNED decoder than for the BASE one as the amount of noise increases. We note also that the overall score distribution (after averaging over all α), shown in Figure 15, exhibits a large spread of fragility scores across SONAR samples.

One factor that can explain such a discrepancy is the text length. Compared to the Auto-Encoding BLEU metric (which drops only by 1–2% for long sentences), fragility is more sensitive to the length of sentences and drops faster for both similarity metrics. This shows that using a max sentence length over 250 can be extremely challenging for SONAR and the LCM model. On the other hand, even if short sentences are on average more robust, splitting a long sentence in the wrong place may result in shorter but more fragile sub-sentences.

Taking a closer look at the 5% most fragile embeddings, we notice that they are very noisy. Typically, they correspond to hyperlinks, references, unique ids, code-switched or numerical entries. These are likely artifacts that the SONAR models were not exposed to during training or where the SONAR tokenizer fails. Fragility can thus be used to filter out hard samples from the training data. We also observe that short but complex technical phrases can be more fragile than common language phrases of similar length.

3 Scaling the model to 7B

This section describes our effort to scale our model to 7B parameters and compare the performance against other approaches such as token-based LLMs on more challenging tasks such as summarization and summarization expansion (detailed in Section 3.1).

Based on the results in Section 2.4 where the two diffusion-based LCMS (ONE-TOWER and TWO-TOWER) outperform the other variants, we decided to scale a diffusion model to 7B parameters. We chose to scale TWO-TOWER given its smaller memory footprint, particularly when processing long contexts with a shallower contextualizer tower

The large 7B TWO-TOWER diffusion LCM has 5 layers in its contextualizer and 14 layers in its denoiser. Its dimension has been extended to $d_{\text{model}}=4096$. Each self-attention layer has 32 attention heads. All other parameters are kept the same as for the 1.6B TWO-TOWER model. The model is

pre-trained on a dataset of 2.3B documents, representing 2.7T tokens and 142.4B concepts/sentences. We pre-trained this model on Meta’s RSC for 124k optimization steps spanning 256 A100 GPUs with a total batch size of 1M concepts. We further extend the context length of this model to cover 2048 concepts instead of the 128 concepts in the ablation experiments. We trained using the AdamW optimizer with $(\beta_1, \beta_2) = (0.9, 0.95)$, $\epsilon = 1e-5$ and a weight decay of 0.1. We use a cosine learning rate schedule, with warm-up of 10,000 steps up to $LR = 3e-4$. To improve training stability we clip gradients at a maximum norm of $g = 10$. We subsequently finetune the 7B TWO-TOWER LCM on publicly available instruction tuning datasets following [Chung et al. \(2024\)](#). Each sample consists of a prompt and an answer and we back-propagate on answer sentences only. Each answer sequence is suffixed with the phrase “End of response.” to teach the model when to stop generating. The finetuning data totals 389M sentences, of which 53M are answers (*i.e.*, targets). For supervised finetuning, we use a cosine learning rate schedule with an initial rate of $LR = 3e-5$ and finetune the model for 7 epochs with a batch size of 262K sentences (prompts and answers combined). We will refer to the pre-trained model as TWO-TOWER-7B and the finetuned model as TWO-TOWER-7B-IT.

3.1 Evaluation Tasks and Data

This section describes the tasks on which we are evaluating and benchmarking our proposed model. We detail datasets, baselines and metrics. For each task, the dataset was processed with the same sentence splitter and SONAR encoder as used in the LCM training.

3.1.1 Metrics

As longform text generation is the main challenge for LCM, our benchmarking is mainly focused on generative tasks, which are notoriously more difficult to evaluate automatically. Therefore, we evaluate them with multiple automatic metrics, chosen to focus on complementary aspects on generation quality. All metrics used in this section are summarized in [Table 8](#).

For summarization and summary expansion (defined below), we report the traditional reference-based ROUGE-L metric ([Lin, 2004](#)). As summarization models have a tendency to copy content from the source or from its own generated prefix, we report two additional word-based metrics. To evaluate how much content is directly copied from the source, we report the proportion of word 3-grams of the source that are present in the output (OVL-3). To evaluate repetitiveness of the generated texts, we report the portion of duplicated word 4-grams in the output (REP-4).

To complement word-based metrics with summarization-focused neural evaluation, we use two metrics introduced by [Clark et al. \(2023\)](#): average probabilities of the SEAHORSE classifiers for Q4 (whether all the information in the summary is fully attributable to the source), denoted as SH-4 in the following and Q5 (whether the summary captures the main ideas of the source), denoted as SH-5.

As a metric of the overall fluency of the generated sentences, we report an average probability that the sentence is linguistically acceptable, as predicted by a classifier trained by [Krishna et al. \(2020\)](#) on the COLA dataset ([Warstadt et al., 2019](#)), further referred to as COLA. To evaluate the local coherence of the generated text, we report the average cosine similarity between each n ’th and $n + 2$ ’th sentence ([Parola et al., 2023](#)).

3.1.2 Summarization

Task and datasets. When considering a relatively long document, a summarization task can be described as the act of generating a much shorter corresponding document that includes the essential

Task	Area	Metric	Description	Reference
Summarization	Target similarity	R-L	ROUGE-L	Lin (2004)
	Source similarity	OVL-3	N-grams overlap (N=3)	
	Grammaticality	REP-4	Portion of duplicated N-grams (N=4)	Welleck et al. (2019)
	Fluency	CoLA	Sentence fluency classifier score	Krishna et al. (2020)
	Attribution	SH-4	Seahorse-Large-Q4 score	Clark et al. (2023)
	Semantic coverage	SH-5	Seahorse-Large-Q5 coverage score	Clark et al. (2023)
Summary	Grammaticality	REP-4	(see above)	Welleck et al. (2019)
Expansion	Fluency	CoLA	(see above)	Krishna et al. (2020)

Table 8 - Summary of automatic metrics used in different tasks in Section 3.1. Order mostly follows paper’s narrative.

information contained in the long document and the same logical structure linking the various pieces of essential information.

Summarization techniques can range from more extractive to more abstractive. Extractive techniques attempt to preserve in the summary the same vocabulary as that found in the long document, thereby shortening the long by removing details and superfluous wording. Abstractive techniques, on the other hand, attempt to produce the summary by rephrasing the essential pieces of information found in the long document. Our work focuses more on abstractive summarization, as such type of summarization cannot be performed without some form of understanding and reasoning.

We use the CNN DAILYMAIL (Hermann et al., 2015) and XSUM (Narayan et al., 2018) datasets. We also report results on the challenging LCFO corpus which takes long documents as input, approx. 5k words (Costa-jussà et al., 2024). The task is to provide abstractive summaries with lengths representing 20%, 10%, and 5% of the input document. Detailed statistics are provided in Table 9.

Dataset	#Docs	#LLAMA2 Tokens			#Sentences		
		Q1	Q2	Q3	Q1	Q2	Q3
CNN DAILYMAIL	11.5k	605/61	892/78	1266/97	10/3	14/4	21/4
XSUM	11.3k	273/25	445/30	735/35	25/1	30/1	35/1
LCFO.5%	249	6559/341	7214/378	7916/418	209/12	295/15	527/18
LCFO.10%	249	6559/654	7214/718	7916/796	209/22	295/27	527/32
LCFO.20%	249	6559/1276	7214/1403	7916/1524	209/41	295/48	527/59

Table 9 - Statistics of the test split of evaluation benchmarks. For each subset we report the number of documents and statistics of document and summary length in terms of sentences and LLAMA2 tokens. Each table cell shows “document/summary” length quartiles.

Baselines. For CNN DAILYMAIL and XSUM, we compare against several baselines of different architectures (encoder-decoder transformer, decoder-only LLMs) that are known to perform well on summarization tasks. For encoder-decoder transformer models, we use T5 (Raffel et al., 2020). For decoder-only LLMs, we choose GEMMA-7B, LLAMA-3.1-8B and MISTRAL-7B-v0.3. We chose the published instruction-tuned models to compare with the LCM with the same training regime, and have a similar size (7B). Note that while T5 has much smaller sizes than the LCM, this is compensated by using models that are fine-tuned explicitly on the target evaluation dataset.

Summarization results. Table 10 contains the results of different baselines and our LCM model for summarization (CNN DAILYMAIL and XSUM). We can notice that the LCM produces competitive ROUGE-L scores when compared to a specifically tuned LLM (T5-3B) and even surpasses the

MODEL	Paradigm	CNN DAILYMAIL					
		R-L(↑)	OVL-3 (↑)	REP-4 (↓)	CoLA (↑)	SH-4 (↑)	SH-5 (↑)
Ground truth	—	100.00	0.170	0.684	0.850	0.683	0.586
T5-3B	SFT	37.56	0.174	0.854	0.946	0.773	0.503
GEMMA-7B-IT	IFT	31.14	0.245	1.032	0.963	0.740	0.560
MISTRAL-7B-v0.3-IT	IFT	36.06	0.200	0.780	0.972	0.780	0.676
LLAMA-3.1-8B-IT	IFT	34.97	0.248	0.928	0.973	0.763	0.692
TWO-TOWER-7B-IT	IFT	36.47	0.177	0.757	0.767	0.723	0.459

MODEL	Paradigm	XSUM					
		R-L(↑)	OVL-3 (↑)	REP-4 (↓)	CoLA (↑)	SH-4 (↑)	SH-5 (↑)
Ground truth	—	100.00	0.108	0.399	0.987	0.352	0.418
T5-3B	—	17.11	0.221	0.671	0.939	0.680	0.450
GEMMA-7B-IT	IFT	18.20	0.177	0.620	0.769	0.546	0.446
MISTRAL-7B-v0.3-IT	IFT	21.22	0.162	0.480	0.922	0.633	0.621
LLAMA-3.1-8B-IT	IFT	20.35	0.186	0.501	0.941	0.687	0.658
TWO-TOWER-7B-IT	IFT	23.71	0.106	0.464	0.683	0.358	0.284

Table 10 - Performance on the CNN DAILYMAIL and XSUM summarization tasks.

instruct-finetuned LLMs. Our model tends to generate more abstractive summaries rather than extractive ones, as shown by the lower OVL-3 scores. The LCM produces fewer repetitions compared to LLMs, and more importantly, the repetition rate is closer to the ground truth one. The LCM generates globally less fluent summaries according to CoLA classifier. However, we can remark that even the human generated ground truth gets a lower score compared to the LLM. A similar behavior is observed for the source attribution (SH-4) and semantic coverage (SH-5). This may be explained by model-based metrics that are more biased towards LLM generated content.

Long-context summarization results. Table 11 presents the results for long-context summarization (LCFO.5%, LCFO.10% and LCFO.20%). This is a challenging task for most of the models. For example, MISTRAL-7B-v0.3-IT seems to be unable to follow the length instruction of the summary—it always generates summaries which length is about 50% of the source. MISTRAL-7B-v0.3-IT also has the highest SH-4 score, *i.e.*, source attribution. The summaries generated by GEMMA-7B-IT tend to be longer than requested, while LLAMA-3.1-8B-IT generates summaries which length is the closest to the requested size.

The LCM has only seen a limited amount of long documents in the pretraining and fine-tuning data. Nevertheless, it performs well for this task. It outperforms MISTRAL-7B-v0.3-IT and GEMMA-7B-IT in the metric ROUGE-L for the 5 and 10% conditions, and is close to GEMMA-7B-IT for the 20% condition. We also observe that the LCM yields high SH-5 scores for all conditions, *i.e.*, the summaries can be attributed to the source.

Finally, we observe that LLAMA-3.1-8B-IT performs substantially better than the other LLMs, according to ROUGE-L, while all LLMs have similar performance on the CNN DAILYMAIL and XSUM summarization tasks. This could be explained by training data contamination for LLAMA-3.1-8B-IT, or by the fact that the other two LLMs struggle to handle the long input context.

METHOD	WR	LCFO.5%					
		R-L(↑)	OVL-3 (↑)	REP-4 (↓)	CoLA (↑)	SH-4 (↑)	SH-5 (↑)
GEMMA-7B-IT	0.107	25.21	0.151	4.711	0.688	0.357	0.174
MISTRAL-7B-v0.3-IT	0.512	21.36	0.532	5.997	0.854	0.656	0.296
LLAMA-3.1-8B-IT	0.076	37.67	0.190	2.767	0.931	0.488	0.314
TWO-TOWER-7B-IT	0.060	26.88	0.162	2.473	0.796	0.628	0.196
		LCFO.10%					
		R-L(↑)	OVL-3 (↑)	REP-4 (↓)	CoLA (↑)	SH-4 (↑)	SH-5 (↑)
GEMMA-7B-IT	0.150	29.25	0.164	6.427	0.667	0.377	0.194
MISTRAL-7B-v0.3-IT	0.549	25.00	0.537	6.289	0.848	0.660	0.306
LLAMA-3.1-8B-IT	0.128	42.85	0.243	3.804	0.907	0.486	0.310
TWO-TOWER-7B-IT	0.089	29.38	0.202	3.00	0.791	0.623	0.183
		LCFO.20%					
		R-L(↑)	OVL-3 (↑)	REP-4 (↓)	CoLA (↑)	SH-4 (↑)	SH-5 (↑)
GEMMA-7B-IT	0.257	33.32	0.201	9.188	0.603	0.425	0.239
MISTRAL-7B-v0.3-IT	0.493	28.82	0.527	5.806	0.858	0.658	0.293
LLAMA-3.1-8B-IT	0.179	46.92	0.272	4.783	0.888	0.485	0.315
TWO-TOWER-7B-IT	0.140	31.74	0.253	3.664	0.779	0.613	0.187

Table 11 - Performance on the long-context summarization task of LCFO. WR is the word count ratio between the generated text and the source document.

4 Large Concept Model Extensions

In this section, we explore several extension of the LARGE CONCEPT MODEL. First, we evaluate the LCM on the new task of summary expansion, *i.e.*, given a summary, create a longer text. We then showcase the good zero-shot generalization performance of the LCM. Finally, we explore an approach to add higher level information beyond sentences.

4.1 Summary Expansion

Task and datasets. When considering a short and concise document that has similar properties to those of a summary (*i.e.*, mainly a stand-alone document that abstracts from details), a summary expansion task can be described as the act of generating a much longer document that preserves the essential elements found in the corresponding short document, as well as the logical structure that connects such elements. As this is a more freely generative task, an additional requirement to be taken into consideration is that of coherence (for example, the detailed information included in one generated sentence should not contradict that included in another sentence). The summary expansion task presented here consists in taking summaries as inputs, from CNN DAILYMAIL and XSUM, and generating a long document. Note that the goal is not to recreate the factual information of the initial document rather than evaluating the capability of the model to extend the input text in a meaningful and fluent way. We use similar baselines and metrics as in the previous section 3.1.2.

Results. Table 12 shows the results of the summary expansion for CNN DAILYMAIL and XSUM. First of all, regarding the word count ratio, we can see different behaviours for the two corpora. For CNN DAILYMAIL, the models tend to generate texts that are 6 times larger than the input. LLAMA-3.1-8B-IT produces even longer outputs (factor 8 instead of 6). But for XSUM, while

CNN DAILYMAIL					
METHOD	WR	R-L(↑)	OVL-3 (↑)	REP-4 (↓)	CoLA (↑)
GEMMA-7B-IT	6.8	35.54	0.801	2.104	0.951
MISTRAL-7B-v0.3-IT	6.4	34.24	0.817	2.063	0.959
LLAMA-3.1-8B-IT	8.5	37.76	0.822	2.582	0.844
TWO-TOWER-7B-IT	6.3	30.85	0.726	2.911	0.474
XSUM					
METHOD	WR	R-L(↑)	OVL-3 (↑)	REP-4 (↓)	CoLA (↑)
GEMMA-7B-IT	19.5	17.89	0.963	10.238	0.116
MISTRAL-7B-v0.3-IT	1.6	29.31	0.893	2.268	0.939
LLAMA-3.1-8B-IT	19.8	28.84	0.915	2.543	0.898
TWO-TOWER-7B-IT	7.1	23.82	0.561	1.542	0.603

Table 12 - Performance on the summary expansion tasks of CNN DAILYMAIL and XSUM, evaluated with the metrics described in Table 8. WR is the word count ratio between the hypothesis and the source summary.

GEMMA-7B-IT and LLAMA-3.1-8B-IT generate very long texts (almost 20 times longer than the prompt), the LCM generates output of approximately the same length ratio as for CNN DAILYMAIL. Only MISTRAL-7B-v0.3-IT fails to generate long outputs for this corpus.

Then, we clearly see a different trend compared to the summarization task. The LLMs get higher ROUGE-L scores compared to the LCM. As mentioned above, the goal of this task is not to recreate the original document. However, the R-L score tells us how much of the content of the full document can be recreated. Contrary to the LLMs, our model tends to generate different sentences compared to the original document from which the summary has been created (with an exception for GEMMA-7B-IT on XSUM). This is expected since our model generates embeddings that are then processed by a decoder trained on a translation task that tends to paraphrase the initial content. However, the CoLA results show that this comes along with lower fluency, especially for CNN DAILYMAIL.

4.2 Zero-shot generalization performance

SONAR is a semantic space that can represent 200 languages. In this paper, all experiments presented so far have been done on English text. In this section, we explore the capability of our proposed LCM approach to process other languages in a zero-shot fashion by leveraging SONAR’s ability to represent multilingual data.

We use the XLSUM (Hasan et al., 2021) corpus, a large scale multilingual abstractive news summarization benchmark covering 45 languages. We score model outputs using the multilingual rouge scoring scripts released with the benchmark.⁵ Note that ROUGE-L scores are heavily dependent on language-specific text tokenization and stemming. Unless provided in the aforementioned scripts, we tokenize the model outputs and references with the default tokenization from Lin (2004). Languages like Korean, Telugu and Tamil can benefit from a more appropriate stemming and tokenization.

We compare the LCM performance with LLAMA-3.1-8B-IT which officially supports eight languages: English, German, French, Italian, Portuguese, Hindi, Spanish, and Thai. According to The Llama3 team (2024), the model has seen many additional languages during pretraining, but was instruction

⁵https://github.com/csebuetnlp/xl-sum/tree/master/multilingual_rouge_scoring

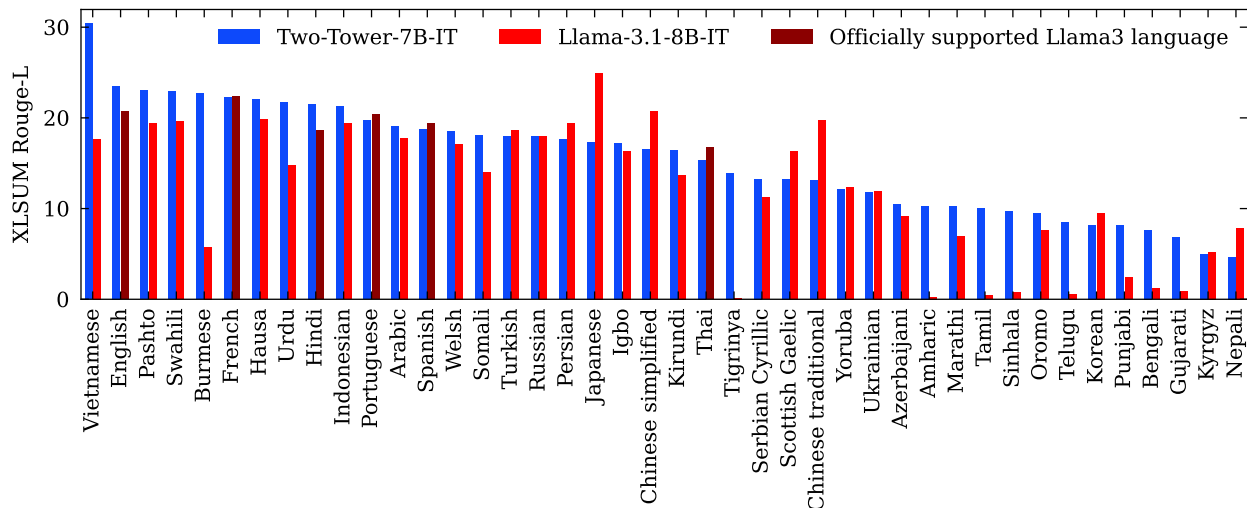


Figure 16 - ROUGE-L scores on XLSUM for LLAMA-3.1-8B-IT and TWO-TOWER-7B-IT.

finetuned on those eight languages only. The LCM, on the other hand, has never seen any language other than English, and we do not use the English XLSUM training data either.

We report ROUGE-L scores for 42 languages in Figure 16. Three languages were excluded since they are currently not supported by SONAR: Pidgin, Serbian in Latin script and Uzbek in Cyrillic script. The LCM substantially outperforms LLAMA-3.1-8B-IT on English (23.5 compared to 20.7 ROUGE-L) and on the average over all six languages officially supported by both models and included in XLSUM (20.2 versus 19.7 ROUGE-L).⁶ We also observe that the LCM generalizes very well to many other languages, in particular low-resource languages like Southern Pashto, Burmese, Hausa or Welsh which all have ROUGE-L scores greater than 20. Other well performing low-resource languages are Somali, Igbo or Kirundi. Finally, the LCM obtains a ROUGE-L score of 30.4 on Vietnamese. Overall, these results highlight the impressive zero-shot generalization performance of the LCM to languages it has never seen.

4.3 Exploring explicit planning

When writing long-form text, it is often practical to first think about how to structure our narrative. Indeed Dijk (1977) states that all texts have an innate macrostructure *i.e.*, a global discourse structure spanning the entirety of the text. In general scientific discourse, one such macrostructure is that of problem-solving (Heffernan and Teufel, 2022), where an author must first motivate and describe the problem, before proposing an appropriate solution.

Composing such a macrostructure is no easy task. Yet in the realm of LLMs, there have been many recent efforts to help guide model generation using similar structures. One such popular approach is that of creating outlines, which provide a high-level overview of the desired narrative (Li et al., 2024). An alternative approach to outlines is creating summaries of future targets, which the model can then expand upon (Sun et al., 2022).

Given the nature of the LCM operating at the concept level, it naturally creates long-form output. Therefore, it is important to ensure that the model is capable of creating coherent generations given the multitudes of possibilities for next concept prediction. In order to address this, we envision an

⁶English, French, Hindi, Portuguese, Spanish and Thai.

explicit capability for planning. Similar to creating summaries, we propose a complementary *planning model* which creates a high-level overview of what should be generated next, given the prior context. The proposed plan could span multiple concepts, such as a paragraph. The LCM is conditioned on this plan, before it then generates the subsequent output sequence.

Operationally the model predicts auto-regressively a sequence of concepts followed by a *break* concept, which represents a natural topic cadence such as a paragraph break. Once the *break* concept is predicted, the large planning model (LPM) generates a plan in order to condition the LCM for prediction of the subsequent sequence. The model then continues generation as usual conditioned on both the prior concepts and the proposed plan. An overview of the approach is shown in Figure 17.

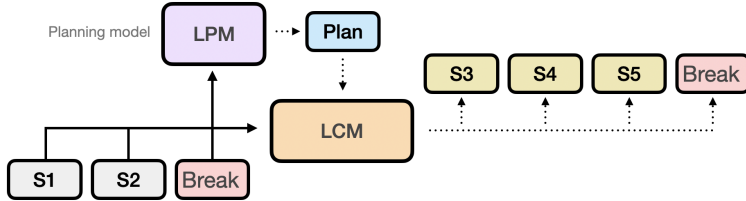


Figure 17 - LCM conditioned on both the prior context and a high-level plan for the next sequence.

Although we envision a separate (but complementary) model for planning, we present here an initial experiment using a simplified single-model approach where the LCM is trained in a multitask setting to also predict both the *break* concepts and plans. Additionally, instead of the idealized plan spanning multiple concepts (such as a paragraph), we use a single concept in order to capture what should come next (i.e. a *plan* concept). We call this simplified multitask approach a Large Planning Concept Model (LPCM).

Methodology

In order to evaluate this single-model approach, we perform an ablation study. As a baseline method, we train a ONE-TOWER LCM (cf. Section 2.3.3) without any visibility to *break* or *plan* concepts. We then subsequently train a LPCM with the same number of parameters as the baseline. Both models were trained on the same data⁷ for the same number of steps.

Data preprocessing. In order to represent *break* concepts, we begin by first segmenting the data into paragraphs. Given that most real world datasets are absent of paragraph structure (or it is not easy to recover), we apply the Segment Any Text (Frohmann et al., 2024) paragraph splitting API⁸. We additionally force each paragraph to be less than 10 sentences, and merge small (e.g. one sentence) consecutive paragraphs together. In order to represent *plan* concepts, we generate synthetic high-level topic description for each preceding segmented paragraph using an existing open-sourced LLM, namely Llama-3.1-8B-IT, which offers a good trade-off between the generated topic quality and the generation speed. The system prompt used to generate these topic descriptions is listed in Appendix C. In total we process approximately 320M paragraphs with topic descriptions, spanning 1.5B segmented concepts (i.e. approximately 30B tokens).

Metrics. We focus on coherence as our main measure of evaluation. Previous ablations (cf. Section 2.4) used the coherence metric introduced by Jwalapuram et al. (2022). However, we explore

⁷Compared to previous experiments, we use a different data mixture more favorable for long-form generation.

⁸<https://github.com/segment-any-text/wtpsplit>

here LLM-as-a-judge as an alternative. Specifically, we use Llama-3.1-8B-IT in order to evaluate the coherence of the generated model outputs, which is prompted to return an overall coherence score between [0, 5]. The prompt used is listed in [Appendix D](#). In order to validate this prompt, we evaluate it against a dataset of human judgements introduced by [Jwalapuram et al. \(2022\)](#), and observed it reported an agreement⁹ with human annotators which improves upon their coherence model. We therefore choose this metric for our evaluation. To be consistent across both model results, we do not include the special *break* or *plan* concepts generated by the LPCM when calculating coherence scores.

Llama-3.1-8B-IT (\uparrow)	
LPCM	2.82 \pm 0.62
Baseline	2.74 \pm 0.70

Table 13 - LPCM ablation coherence score results.

Results

We provide the results of our ablation experiment in [Table 13](#). Results are reported over a held-out subset of COSMOPEDIA ([Ben Allal et al., 2024](#)) following instruction fine-tuning, similar to previous ablations (cf. [Section 2.4](#)). We observe that the LPCM achieves significantly¹⁰ higher coherence scores (significance was measured using a paired t-test) than the baseline ONE-TOWER LCM. This finding suggests that the LPCM is capable of producing significantly more coherent outputs than the LCM as a result of the additional structure coming from predicted plan concepts, helping the LPCM produce a more coherent narrative, which is essential for the objective of generating long-form output.

5 Related work

5.1 Sentence representations

Multilingual sentence representations Learning effective sentence embeddings has been a well studied subject in recent years. Significant progress has been made in this field, largely due to the capabilities of transformer-based language models that by learning contextual representations for individual tokens ([Devlin et al., 2018](#); [Conneau et al., 2020](#)), are able to effectively learn the semantics of language. However, such models are not optimal to create sentence representations.

Following approaches built upon these initial works, and aimed at learning general sentence representations, leveraging dual encoder architectures ([Guo et al., 2018](#); [Reimers and Gurevych, 2019](#); [Ni et al., 2021](#)). These architectures encode the source and target into a common embedding space, and use a distance-based metric to create an alignment loss that approximate semantically identical sentences. Such architectures have been extended to leverage multilingual data to create general, aligned embedding spaces across languages ([Feng et al., 2020](#); [Janeiro et al., 2024](#); [Sturua et al., 2024](#)). Initial approaches leveraged the contrastive loss to align translations across languages ([Feng et al., 2020](#); [Yang et al., 2019](#)), using only translation data to train. Other architectural changes, namely using token-level objectives combined with the sentence level objectives, have proven useful to improve the quality of multilingual sentence representations based on translation data only ([Li](#)

⁹Krippendorff’s $\alpha = 0.48$

¹⁰Significance observed at the 99.9% level.

et al., 2023; Janeiro et al., 2024). Recent approaches explore using data from other tasks, besides translation data, to increase the generality of the sentence representations (Wang et al., 2024b; Mohr et al., 2024). Other approaches change their embeddings per task, either with task-specific prompts (Wang et al., 2024b; Su et al., 2022; Lee et al., 2024b) or with task-specific parameters (Sturua et al., 2024).

Another successful line of work to create general purpose, multilingual, sentence representations is to leverage the translation objective. LASER (Artetxe and Schwenk, 2019), and SONAR (Duquenne et al., 2023b) leverage an encoder-decoder architecture, with a fixed-size sentence representation between the encoder and the decoder, trained with a translation objective. SONAR is initialized from the NLLB-200 model (NLLB Team et al., 2022), and covers 200 languages, making it one of the open-source models with the widest language coverage. SONAR also provides open-source speech encoders aligned to their sentence encoders for 73 languages (Seamless Communication et al., 2023a), aligned through a teacher-student approach. SONAR has been used as the basis for several works (Seamless Communication et al., 2023a,b; Chen et al., 2023a), and its speech decoders have been extended to keep the expressiveness of the original speech (Duquenne et al., 2023a).

Joint speech/text sentence representations There has been a large body of research on unsupervised representation learning for monolingual (Baeviski et al., 2020) and multilingual speech (Babu et al., 2022), with recently w2v-bert (Chung et al., 2021) that combines contrastive learning and masked language modeling to learn self-supervised representations from speech. Other works explored multilingual and multimodal (speech/text) pre-training methods, including mSLAM (Bapna et al., 2022). Finally, Duquenne et al. (2021), followed by Khurana et al. (2022), introduced multilingual and multimodal sentence embeddings, extending a pre-existing multilingual text sentence embedding space to the speech modality with a distillation approach. Duquenne et al. (2022, 2023c) also showed that it is possible to efficiently decode multilingual speech sentence embeddings with decoders trained on text sentence embeddings into different languages, to perform zero-shot speech translation.

LLM based sentence representations Several text representation methods have been proposed which are based on existing LLMs. Wang et al. (2024a) proposed extracting text embeddings from the last token of LLMs fine-tuned with instructions on contrastive data. Lee et al. (2024a) improved text embedding capabilities of fine-tuned LLMs by removing the causal attention mask and applying extra nonlinear layers before pooling the token embeddings. Embeddings as a service are supported by some commercial LLM providers, for example, MISTRAL-EMBED.¹¹ Such embeddings proved competitive on retrieval benchmarks; however, to the best of our knowledge, their applicability to reconstructing the texts back from the embedding space has not been demonstrated.

5.2 Multilingual LLMs

Most of the leading LLMs have been trained on texts in several languages. Table 1 summarizes the coverage of several of them. Nevertheless, the pretraining data of these LLMs seems to be mainly English texts. For example, The Llama3 team (2024) mentions that pretraining data contains significantly more English texts, requiring continued pre-training with multilingual data, out of which 34.6% is translated reasoning data.

There are also several efforts to train LLMs optimized on specific languages, e.g. LEOLM for German,¹² FUANO for Italian (Bacciu et al., 2024), ALLAM for Arabic (Bari et al., 2024), and several

¹¹<https://docs.mistral.ai/capabilities/embeddings/>

¹²<https://laion.ai/blog/leo-lm/>

models for Chinese: ERNIBOT,¹³ TONGYI QIANWEN,¹⁴ or CHATGLM (Team GLM et al., 2024). Some adaptations of LLMs to a massive number of languages also exist. LOLA (Srivastava et al., 2024) is a recent mixture-of-experts LLM supporting 160 languages, MALA-500 (Ji et al., 2024) adapts LLaMA2 to 546 languages. However, such models typically face a trade-off between language coverage and other capabilities. For example, the Aya model (Üstün et al., 2024), following instructions in 101 languages, was superseded by Aya-23 (Aryabumi et al., 2024) that exchanged some breadth for depth, focusing on 23 languages only. The LCM architecture, combining a language-agnostic model for knowledge and reasoning with potentially language-specialized encoders and decoders, is expected to exhibit this trade-off to a lesser extent.

5.3 Alternative LLM architectures

Predicting the next state in the embedding space is a core idea of the Joint Embedding Predictive Architecture (JEPA) proposed by LeCun (2022). This idea has been implemented for images (I-JEPA by Assran et al. (2024)) and video (V-JEPA by Bardes et al. (2024)) as a self-supervised approach to learning representations. For language, equivalent models have not yet been explored.

Sentence embeddings for language modeling. For text completion, Ippolito et al. (2020) proposed a sentence-level language model operating by choosing the next sentence from a finite set of candidates. Their model demonstrated success in selecting appropriate continuations for short stories, but it has not been scaled to longer inputs or to fully generative outputs. Golestani et al. (2021) studied a similar problem in the even more restrictive sentence ordered setting, but with a more thorough study of architectural choices. The INSET architecture (Huang et al., 2020) solves the sentence infilling task by combining a denoising autoencoder that encodes sentences into fixed-size vectors and decodes them back and a bidirectional transformer that predicts the embedding of a missing sentence.

Marfurt and Henderson (2021) and Cornille et al. (2024) used predicted next sentence embeddings in a fully generative setting, for summarization and generic language modeling, respectively. However, their architectures considered sentence-level connections only as an addition to the token-level connections across sentences, not as their replacement.

In a recent work of An et al. (2024), the SentenceVAE architecture performs language modeling on the sentence level using a sentence encoder to prepare the inputs and a sentence decoder to produce the outputs. However, its input and output embedding spaces are not tied, so the inference is only possible by decoding each predicted sentence into text and then re-encoding it for adding it to the context.

Language modeling with diffusion. A series of more recent works tried adapting diffusion modeling, originally developed for continuous data, to the discrete text domain. The PLANNER architecture (Zhang et al., 2023) consists of a variational autoencoder for paragraphs and a diffusion model trained to predict latent autoencoder representations conditional on the textual context or on the class label. Lovelace et al. (2024) augmented a decoder-only language model with an encoded semantic proposal of the continuation text, with an easily guidable diffusion model predicting the embedding of the next proposal. A TEncDM model (Shabalin et al., 2024) performs diffusion in the space of contextual token embeddings which are then decoded non-autoregressively.

¹³<http://research.baidu.com/Blog/index-view?id=183>

¹⁴<https://www.alibabacloud.com/en/solutions/generative-ai>

Some applications of diffusion to sequence modeling have targeted the planning capabilities of the sequence models. Semformer (Yin et al., 2024) proposed training transformers language models to plan several steps ahead by including special planning tokens, the representations of which are trained to be informative about the future tokens. Ye et al. (2024) applied discrete diffusion to language models as an alternative to autoregressive generation, more suitable for tasks that require multi-step planning. Ubukata et al. (2024) give an overview of applications of diffusion for planning tasks, but most of them are not concerned with the language domain.

Overall, while many of the previous works used hidden representations for language modeling or related tasks, all of them either relied on token-level inputs or outputs, or were not intended for generating texts of arbitrary length. The LCM seems to be the first fully generative language model implemented fully in a highly semantic, reconstructable sentence representation space.

6 Limitations

In this section we discuss the possible limitations of the presented Large Concept Modeling approach.

Choice of the embedding space. The choice and design of the embedding space plays a crucial role in the LCM modeling approach.

- The SONAR embedding space was chosen for its good multilingual and multimodal representations, as well as the availability of a massively multilingual decoder, which achieves excellent results in both translation and auto-encoding. However, the SONAR model was trained on very specific training data, namely bitext machine translation data containing rather short sentences. This has several consequences:
 1. SONAR is trained to sustain a local geometry (sentences with very similar meanings are geometrically close) with no special guarantees for sentences that are only loosely related. Yet, predicting next sentences distribution requires the space to operate well globally.
 2. SONAR auto-encodes surprisingly well texts containing links, references, or merely numbers or code data. Yet, such texts tend to be fragile, highlighting a distribution mismatch between the SONAR training data and commonly used LLM pre-training text corpora. Therefore, the accurate prediction of the sentences containing such a content (non-negligible in LCM pre-training data) will be hard for any LCM SONAR based model. For instance, the factuality of fragile generated sentences may easily be compromised.
- Using a frozen encoder represents some interesting trade-offs. Any frozen encoder which is learned in a different data context, and with no a-priori strong connection to LCM modeling, may be suboptimal compared to encoders that are learned in an end-to-end fashion (with the loss coming from the decoder). At the same time, learning an encoder within end-to-end training can be challenging and the resulting space is not guaranteed to result in good semantic representations shared across languages and modalities.

Training the concept representation and the LCM end-to-end would also be less data and compute efficient since all modeling data should be multilingual and -modal, bearing the risk of modality competition.

Concept granularity

- In this work, the definition of concepts is interpreted at sentence level. However, the manifold of possible next sentences is very wide, attributing a proper probability to each of such sentences is much harder (even with a modeling within the latent space) than to the discrete set of tokens.
- In NLP, we encounter sentences of variable length. Combinatorial complexity of possible next sentences grows exponentially with the maximum character length. The choice of granularity for LCM is not trivial as long sentences (>120 characters) could reasonably be considered as several concepts. However, any finer splitting of such sentences does not necessarily separate well these concepts. This shows the limitation of a fixed size embedding representation for one sentence. Text splitting (such as sentence splitting) or one-to-many mapping of a sentence into several embeddings is a major future direction of research.
- Each document in a training corpus typically contains a sequence of unique sentences or a little number of repetitions. This data sparsity effect manifests as well at large corpora level: the large majority of sentences are merely unique. In principle, this issue can be addressed with higher-level semantic embedding representations. These higher-level representations come with trade-off between requirement of lossless data encoding (think of named-entities or numbers, critical in many language modeling tasks) and good level of abstraction to enable reasoning capabilities. Compared to a monolingual auto-encoder which would simply compress input sentences, SONAR offers semantic representations with good auto-encoding quality but still certainly sacrificing generalization capacities.
- This generalization issue can be partially mitigated by splitting or encoding input text as new conceptual units which are more commonly shared across source documents. This is in the spirit of stemming or lemmatization techniques studied in NLP for words. That being said, building such conceptual units that are also language and modality agnostic is a challenging task. Such shared multilingual and multimodal conceptual units are also key for generalization across languages and across modalities. To maximize cross-lingual and cross-modal transfers, LARGE CONCEPT MODELS should be exposed to a richer variety of multilingual and multi-modal data.

Continuous versus discrete

- Diffusion modeling has proven to be very efficient in generative modeling of continuous data like images or speech. As previously stated, sentences in the SONAR space, despite being represented as continuous vectors, remain discrete combinatorial objects. This makes diffusion modeling struggle on the text modality (either at word or sentence embedding level).
- The contrastive nature of cross-entropy loss based on softmax outputs which is used for next token prediction plays a critical role for many downstream tasks where higher accuracy is required (e.g. MCQ tasks, code or math generation). On the opposite, continuous diffusion modeling does not allow to integrate such a contrastive objective.
- The QUANT-LCM could be a way to address the discrete nature of text while modeling on coarse-to-fine semantic units shared across languages and modalities. The limited performance of the QUANT-LCM approaches presented in this paper may be explained by the fact that SONAR space was not trained to be efficiently quantizable, yielding a significant number of codebooks and a large amount of units per codebook. Therefore, the current SONAR quantization suffers from the exponentially increasing number of RVQ units combinations which does not solve the data sparsity/uniqueness issue discussed earlier. This indicates once again the importance of developing a new representation space, either continuous or discrete, for the LARGE CONCEPT MODEL.

7 Acknowledgments

We would like to thank Robbie Adkins, Can Balioglu, Joy Chen, Pascale Fung, Jason Holland, Amita Kamath, Justine Kao, Sagar Miglani, Alice Rakotoarison, Abhilasha Sancheti, Arjang Talattof, Ellen Tan, Carleigh Wood, Shireen Yates, Bokai Yu and Luke Zettlemoyer for comments and suggestions on this work, as well helping us to improve this paper.

8 Conclusion and Future Work

Current best practice for large scale language modeling is to operate at the token level, i.e. to learn to predict the next tokens given a sequence of preceding tokens. There is a large body of research on improvements of LLMs, but most works concentrate on incremental changes and do not question the main underlying architecture. In this paper, we have proposed a new architecture, named a LARGE CONCEPT MODEL (LCM), which substantially differs from current LLMs in two aspects: 1) all modeling is performed in a high-dimensional embedding space instead of on a discrete token representation; and 2) modeling is not instantiated in a particular language or modality, but at a higher semantic and abstract level. We have named the general form of this representation a “*concept*”.

In this paper, to verify the feasibility of the high-level idea, we have assumed that a concept corresponds to a sentence in the text domain, or an equivalent speech segment, and that the embeddings are obtained by the freely available SONAR sentence encoder (Duquenne et al., 2023b). With respect to the specific architecture of the LCM, we have first shown that directly minimizing the MSE loss in the embedding space does not yield good results. We then explored several architectures based on a diffusion process: the ONE-TOWER and TWO-TOWER LCM, as well as a QUANT-LCM which uses quantization of SONAR representations and then modeling on these discrete units. These ablation experiments were performed with models with 1.6B parameters and focused on the generative task of continuing a sequence of sentences. We have then scaled our models to a size of 7B parameters and instruction-finetuned them on several summarization and summary expansion tasks. We provide a detailed comparison to other public models of the same size, namely GEMMA, MISTRAL and LLAMA.

By design, a LCM exhibits strong zero-shot generalization performance. In this paper, we trained models on English texts only, and applied them to text in other languages, without any additional training data, neither aligned nor unlabeled. The LCM outperforms LLAMA-3.1-8B-IT on English and on the average over foreign languages officially supported by the LLM. The LCM itself could also be trained on multilingual- and model data to acquire knowledge from these sources. We will explore this in future versions of the LCM. In short, all languages and modalities are first class citizens and handled equally at all stages of a LCM.

We have observed that next sentence prediction is substantially more challenging than next token prediction. First, given that we operate in an embedding space and at a higher semantic level, the number of possible sentences is virtually unlimited, while token vocabularies are usually in the range of 100k. Second, even given a long context, there is unavoidably more ambiguity in choosing the next sentence than the next token. And third, the usual softmax output layer over the fixed size token vocabulary provides a normalized probability distribution over all possible token continuations. Theoretically, a diffusion process should be able to learn a probability distribution over an output embedding space, but our current experimental evidence indicates that more research is needed to take full advantage of the properties of LARGE CONCEPT MODELS. As an example, the ability to sample multiple embeddings and associate a score would enable beam search to find the best sequence of sentences. Finally, small modeling errors could yield predictions in the embedding space

which do not correspond to valid sentences, i.e. that cannot be decoded into a syntactically and semantically correct sentence. We will work on alternative concept embeddings to SONAR which would be better suited to the next sentence prediction task, and would improve modeling approaches in that concept embedding space.

We see the models and results discussed in this paper as a step towards increasing scientific diversity and a move away from current best practice in large scale language modeling. We acknowledge that there is still a long path to reach the performance of current flagship LLMs. This will require of course further improving the core architecture, but also careful data selection and curation, extensive ablations, optimized and diverse instruction fine-tuning, and finally, scaling to models with more than 70B parameters.

We open-source the full training code of all our LCM variants, together with a set of supporting scripts,¹⁵ to make it easy for other teams to train LCM models. By these means, we hope to foster research on alternative LLMs and contribute to advance the field of machine intelligence.

References

- A. Aghajanyan, L. Yu, A. Conneau, W.-N. Hsu, K. Hambarzumyan, S. Zhang, S. Roller, N. Goyal, O. Levy, and L. Zettlemoyer. Scaling laws for generative mixed-modal language models. In *International Conference on Machine Learning*, pages 265–279. PMLR, 2023.
- E. Almazrouei, H. Alobeidli, A. Alshamsi, A. Cappelli, R. Cojocaru, M. Debbah, Étienne Goffinet, D. Hesslow, J. Launay, Q. Malartic, D. Mazzotta, B. Noune, B. Pannier, and G. Penedo. The Falcon series of open language models. *ArXiv*, abs/2311.16867, 2023. URL <https://arxiv.org/pdf/2311.16867>.
- H. An, Y. Chen, Z. Sun, and X. Li. SentenceVAE: Enable next-sentence prediction for large language models with faster speed, higher accuracy and longer context. *arXiv preprint arXiv:2408.00655*, 2024.
- Anthropic. The Claude 3 model family: Opus, sonnet, haiku, 2024. URL https://www-cdn.anthropic.com/de8ba9b01c9ab7cbabf5c33b80b7bbc618857627/Model_Card_Claude_3.pdf.
- M. Artetxe and H. Schwenk. Massively multilingual sentence embeddings for zero-shot cross-lingual transfer and beyond. *TACL*, pages 597–610, 2019.
- V. Aryabumi, J. Dang, D. Talupuru, S. Dash, D. Cairuz, H. Lin, B. Venkitesh, M. Smith, K. Marchisio, S. Ruder, et al. Aya 23: Open weight releases to further multilingual progress. *arXiv preprint arXiv:2405.15032*, 2024.
- M. Assran, Q. Duval, I. Misra, P. Bojanowski, P. Vincent, M. Rabbat, Y. LeCun, and N. Ballas. Self-supervised learning from images with a joint-embedding predictive architecture. *ArXiv*, abs/2301.08243, 2024. URL <https://arxiv.org/pdf/2301.08243>.
- A. Babu, C. Wang, A. Tjandra, K. Lakhotia, Q. Xu, N. Goyal, K. Singh, P. von Platen, Y. Saraf, J. Pino, A. Baevski, A. Conneau, and M. Auli. XLS-R: Self-supervised Cross-lingual Speech Representation Learning at Scale. In *Proc. Interspeech 2022*, pages 2278–2282, 2022. doi: 10.21437/Interspeech.2022-143.
- A. Bacciu, G. Trappolini, A. Santilli, E. Rodolà, and F. Silvestri. Fauno: The italian large language model that will leave you senza parole! *ArXiv*, abs/2306.14457, 2024. URL <https://arxiv.org/pdf/2306.14457>.
- A. Baevski, Y. Zhou, A. Mohamed, and M. Auli. wav2vec 2.0: A framework for self-supervised learning of speech representations. *NeurIPS*, 33:12449–12460, 2020.
- C. Balioglu. fairseq2, 2023. URL <http://github.com/facebookresearch/fairseq2>.
- A. Bapna, C. Cherry, Y. Zhang, Y. Jia, M. Johnson, Y. Cheng, S. Khanuja, J. Riesa, and A. Conneau. mslam: Massively multilingual joint pre-training for speech and text. *arXiv preprint arXiv:2202.01374*, 2022.

¹⁵https://github.com/facebookresearch/large_concept_model

- A. Bardes, Q. Garrido, J. Ponce, X. Chen, M. Rabbat, Y. LeCun, M. Assran, and N. Ballas. Revisiting feature prediction for learning visual representations from video. *ArXiv*, abs/2404.08471, 2024. URL <https://arxiv.org/pdf/2404.08471>.
- M. S. Bari, Y. Alnumay, N. A. Alzahrani, N. M. Alotaibi, H. A. Alyahya, S. AlRashed, F. A. Mirza, S. Z. Alsubaie, H. A. Alahmed, G. Alabduljabbar, R. Alkhatran, Y. Almushayqih, R. Alnajim, S. Alsubaihi, M. A. Mansour, M. Alrubaian, A. Alammari, Z. Alawami, A. Al-Thubaity, A. Abdelali, J. Kuriakose, A. Abujabal, N. Al-Twairish, A. Alowisheq, and H. Khan. ALLaM: Large language models for arabic and english. *ArXiv*, abs/2407.15390, 2024. URL <https://arxiv.org/pdf/2407.15390>.
- L. Ben Allal, A. Lozhkov, G. Penedo, T. Wolf, and L. von Werra. Cosmopedia, 2024. URL <https://huggingface.co/datasets/HuggingFaceTB/cosmopedia>.
- J. Betker, G. Goh, L. Jing, TimBrooks, J. Wang, L. Li, LongOuyang, JuntangZhuang, JoyceLee, YufeiGuo, WesamManassra, PrafullaDhariwal, CaseyChu, YunxinJiao, and A. Ramesh. Improving image generation with better captions, 2023. URL <https://api.semanticscholar.org/CorpusID:264403242>.
- BigScience Workshop. BLOOM: a 176b-parameter open-access multilingual language model. *ArXiv*, abs/2211.05100, 2023. URL <https://arxiv.org/pdf/2211.05100>.
- Chameleon team. Chameleon: Mixed-modal early-fusion foundation models. *ArXiv*, abs/2405.09818, 2024. URL <https://arxiv.org/pdf/2405.09818>.
- M. Chen, P.-A. Duquenne, P. Andrews, J. Kao, A. Mourachko, H. Schwenk, and M. R. Costa-jussà. BLASER: A text-free speech-to-speech translation evaluation metric. In *ACL*, pages 9064–9079, 2023a. URL <https://aclanthology.org/2023.acl-long.504>.
- M. Chen, K. Heffernan, O. Çelebi, A. Mourachko, and H. Schwenk. xSIM++: An improved proxy to bitext mining performance for low-resource languages. In *ACL*, pages 101–109, 2023b. URL <https://aclanthology.org/2023.acl-short.10>.
- R. Child, S. Gray, A. Radford, and I. Sutskever. Generating long sequences with sparse transformers. *arXiv preprint arXiv:1904.10509*, 2019.
- H. W. Chung, L. Hou, S. Longpre, B. Zoph, Y. Tay, W. Fedus, Y. Li, X. Wang, M. Dehghani, S. Brahma, et al. Scaling instruction-finetuned language models. *Journal of Machine Learning Research*, 25(70):1–53, 2024.
- Y.-A. Chung, Y. Zhang, W. Han, C.-C. Chiu, J. Qin, R. Pang, and Y. Wu. W2v-bert: Combining contrastive learning and masked language modeling for self-supervised speech pre-training. In *2021 IEEE Automatic Speech Recognition and Understanding Workshop (ASRU)*, pages 244–250. IEEE, 2021.
- E. Clark, S. Rijhwani, S. Gehrmann, J. Maynez, R. Aharoni, V. Nikolaev, T. Sellam, A. Siddhant, D. Das, and A. Parikh. Seahorse: A multilingual, multifaceted dataset for summarization evaluation. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 9397–9413, 2023.
- A. Conneau, K. Khandelwal, N. Goyal, V. Chaudhary, G. Wenzek, F. Guzmán, E. Grave, M. Ott, L. Zettlemoyer, and V. Stoyanov. Unsupervised cross-lingual representation learning at scale. In *ACL*, 2020.
- N. Cornille, M.-F. Moens, and F. Mai. Learning to plan for language modeling from unlabeled data. *arXiv preprint arXiv:2404.00614*, 2024.
- M. R. Costa-jussà, P. Andrews, M. C. Megliogli, J. Chen, J. Chuang, D. Dale, C. Ropers, A. Mourachko, E. Sánchez, H. Schwenk, T. Tran, A. Turkatenco, and C. Wood. LCFO: Long context and long form output dataset and benchmarking. *ArXiv*, 2024.
- J. Devlin, M.-W. Chang, K. Lee, and K. Toutanova. BERT: Pre-training of deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805*, 2018.
- P. Dhariwal and A. Nichol. Diffusion models beat gans on image synthesis. *Advances in neural information processing systems*, 34:8780–8794, 2021.

- V. Dijk. *Text and Context: Explorations in the Semantics and Pragmatics of Discourse*. Longman, 1977.
- M. Douze, A. Guzhva, C. Deng, J. Johnson, G. Szilvasy, P.-E. Mazaré, M. Lomeli, L. Hosseini, and H. Jégou. The Faiss library. *ArXiv*, abs/2401.08281, 2024. URL <https://arxiv.org/pdf/2401.08281>.
- P.-A. Duquenne, H. Gong, and H. Schwenk. Multimodal and multilingual embeddings for large-scale speech mining. In *NeurIPS*, volume 34, pages 15748–15761, 2021. URL <https://proceedings.neurips.cc/paper/2021/file/8466f9ace6a9acbe71f75762ffc890f1-Paper.pdf>.
- P.-A. Duquenne, H. Gong, B. Sagot, and H. Schwenk. T-modules: Translation modules for zero-shot cross-modal machine translation. In *EMNLP*, pages 5794–5806, 2022. URL <https://aclanthology.org/2022.emnlp-main.391.pdf>.
- P.-A. Duquenne, K. Heffernan, A. Mourachko, B. Sagot, and H. Schwenk. Sonar expressive: Zero-shot expressive speech-to-speech translation, 2023a.
- P.-A. Duquenne, H. Schwenk, and B. Sagot. SONAR: sentence-level multimodal and language-agnostic representations, 2023b. URL <https://arxiv.org/abs/2308.11466>.
- P.-A. Duquenne, H. Schwenk, and B. Sagot. Modular speech-to-text translation for zero-shot cross-modal transfer. In *Interspeech*, 2023c.
- F. Feng, Y. Yang, D. Cer, N. Arivazhagan, and W. Wang. Language-agnostic BERT sentence embedding. *arXiv preprint arXiv:2007.01852*, 2020.
- M. Frohmann, I. Sterner, I. Vulić, B. Minixhofer, and M. Schedl. Segment Any Text: A universal approach for robust, efficient and adaptable sentence segmentation. In *EMNLP*, pages 11908–11941, 2024. URL <https://aclanthology.org/2024.emnlp-main.665>.
- O. Gafni, A. Polyak, O. Ashual, S. Sheynin, D. Parikh, and Y. Taigman. Make-a-scene: Scene-based text-to-image generation with human priors. In *European Conference on Computer Vision*, pages 89–106. Springer, 2022.
- Gemini Team Google. Gemini 1.5 unlocking multimodal understanding across millions of tokens of conte. *ArXiv*, abs/2403.05530, 2024. URL <https://arxiv.org/pdf/2403.05530>.
- M. Golestani, S. Z. Razavi, Z. Borhanifard, F. Tahmasebian, and H. Faili. Using BERT encoding and sentence-level language model for sentence ordering. In *International Conference on Text, Speech, and Dialogue*, pages 318–330. Springer, 2021.
- P. Goyal. Accurate, large minibatch sg d: training imagenet in 1 hour. *arXiv preprint arXiv:1706.02677*, 2017.
- M. Guo, Q. Shen, Y. Yang, H. Ge, D. Cer, G. H. Abrego, K. Stevens, N. Constant, Y.-H. Sung, B. Strope, et al. Effective parallel corpus mining using bilingual sentence embeddings. *arXiv preprint arXiv:1807.11906*, 2018.
- T. Hang, S. Gu, C. Li, J. Bao, D. Chen, H. Hu, X. Geng, and B. Guo. Efficient diffusion training via min-snr weighting strategy. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 7441–7451, 2023.
- T. Hasan, A. Bhattacharjee, M. S. Islam, K. Mubasshir, Y.-F. Li, Y.-B. Kang, M. S. Rahman, and R. Shahriyar. XL-sum: Large-scale multilingual abstractive summarization for 44 languages. In *Findings of the Association for Computational Linguistics: ACL-IJCNLP 2021*, pages 4693–4703, Online, Aug. 2021. Association for Computational Linguistics. URL <https://aclanthology.org/2021.findings-acl.413>.
- K. Heffernan and S. Teufel. Problem-solving recognition in scientific text. In *Proceedings of the Thirteenth Language Resources and Evaluation Conference*, pages 6045–6058, Marseille, France, June 2022. European Language Resources Association. URL <https://aclanthology.org/2022.lrec-1.650>.
- K. M. Hermann, T. Kocisky, E. Grefenstette, L. Espeholt, W. Kay, M. Suleyman, and P. Blunsom. Teaching machines to read and comprehend. *Advances in neural information processing systems*, 28, 2015.

- J. Ho and T. Salimans. Classifier-free diffusion guidance. *arXiv preprint arXiv:2207.12598*, 2022.
- J. Ho, A. Jain, and P. Abbeel. Denoising diffusion probabilistic models. *CoRR*, abs/2006.11239, 2020. URL <https://arxiv.org/abs/2006.11239>.
- J. Ho, W. Chan, C. Saharia, J. Whang, R. Gao, A. Gritsenko, D. P. Kingma, B. Poole, M. Norouzi, D. J. Fleet, and T. Salimans. Imagen video: High definition video generation with diffusion models. *ArXiv*, abs/2210.02303, 2022. URL <https://arxiv.org/abs/2210.02303>.
- M. Honnibal, I. Montani, S. Van Landeghem, and A. Boyd. spaCy: Industrial-strength Natural Language Processing in Python, 2020.
- Y. Huang, Y. Zhang, O. Elachqar, and Y. Cheng. INSET: Sentence infilling with inter-sentential transformer. In *ACL*, pages 2502–2515, 2020. URL <https://aclanthology.org/2020.acl-main.226.pdf>.
- D. Ippolito, D. Grangier, D. Eck, and C. Callison-Burch. Toward better storylines with sentence-level language models. *ArXiv*, abs/2005.05255, 2020. URL <https://arxiv.org/pdf/2005.05255>.
- J. M. Janeiro, B. Piwowarski, P. Gallinari, and L. Barrault. Mexma: Token-level objectives improve sentence representations, 2024. URL <https://arxiv.org/abs/2409.12737>.
- S. Ji, Z. Li, I. Paul, J. Paavola, P. Lin, P. Chen, D. O’Brien, H. Luo, H. Schütze, J. Tiedemann, et al. Emma-500: Enhancing massively multilingual adaptation of large language models. *arXiv preprint arXiv:2409.17892*, 2024.
- A. Q. Jiang, A. Sablayrolles, A. Roux, A. Mensch, B. Savary, C. Bamford, D. S. Chaplot, D. de las Casas, E. B. Hanna, F. Bressand, G. Lengyel, G. Bour, G. Lample, L. R. Lavaud, L. Saulnier, M.-A. Lachaux, P. Stock, S. Subramanian, S. Yang, S. Antoniak, T. L. Scao, T. Gervet, T. Lavril, T. Wang, T. Lacroix, and W. E. Sayed. Mixtral. *ArXiv*, abs/2401.04088, 2024. URL <https://arxiv.org/pdf/2401.04088>.
- P. Jwalapuram, S. Joty, and X. Lin. Rethinking self-supervision objectives for generalizable coherence modeling. In *ACL*, pages 6044–6059, 2022. URL <https://aclanthology.org/2022.acl-long.418>.
- T. Karras, M. Aittala, T. Aila, and S. Laine. Elucidating the design space of diffusion-based generative models. *Advances in neural information processing systems*, 35:26565–26577, 2022.
- S. Khurana, A. Laurent, and J. Glass. Samu-xlsr: Semantically-aligned multimodal utterance-level cross-lingual speech representation. *arXiv preprint arXiv:2205.08180*, 2022.
- D. Kingma and R. Gao. Understanding diffusion objectives as the elbo with simple data augmentation. *Advances in Neural Information Processing Systems*, 36, 2024.
- N. Kitaev, L. Kaiser, and A. Levskaya. Reformer: The efficient transformer. *arXiv preprint arXiv:2001.04451*, 2020.
- K. Krishna, J. Wieting, and M. Iyyer. Reformulating unsupervised style transfer as paraphrase generation. In *Empirical Methods in Natural Language Processing*, 2020.
- Y. LeCun. A path towards autonomous machine intelligence, 2022. URL <https://openreview.net/pdf?id=BZ5a1r-kVsf>.
- C. Lee, R. Roy, M. Xu, J. Raiman, M. Shoeybi, B. Catanzaro, and W. Ping. NV-Embed: Improved techniques for training LLMs as generalist embedding models. *arXiv preprint arXiv:2405.17428*, 2024a.
- D. Lee, C. Kim, S. Kim, M. Cho, and W.-S. Han. Autoregressive image generation using residual quantization. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 11523–11532, 2022.
- J. Lee, Z. Dai, X. Ren, B. Chen, D. Cer, J. R. Cole, K. Hui, M. Boratko, R. Kapadia, W. Ding, Y. Luan, S. M. K. Duddu, G. H. Abrego, W. Shi, N. Gupta, A. Kusupati, P. Jain, S. R. Jonnalagadda, M.-W. Chang, and I. Naim. Gecko: Versatile text embeddings distilled from large language models, 2024b. URL <https://arxiv.org/abs/2403.20327>.

- K. Lee and S. Sengupta. Introducing the ai research supercluster — meta’s cutting-edge ai supercomputer for ai research, 2022. URL <https://ai.facebook.com/blog/ai-rsc/>.
- Y. Li, Q. Chen, W. Yan, W. Wang, Q. Zhang, and H. Sundaram. Advancing precise outline-conditioned text generation with task duality and explicit outline control, 2024. URL <https://arxiv.org/abs/2305.14459>.
- Z. Li, S. Huang, Z. Zhang, Z.-H. Deng, Q. Lou, H. Huang, J. Jiao, F. Wei, W. Deng, and Q. Zhang. Dual-alignment pre-training for cross-lingual sentence embedding. In A. Rogers, J. Boyd-Graber, and N. Okazaki, editors, *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 3466–3478, Toronto, Canada, July 2023. Association for Computational Linguistics. doi: 10.18653/v1/2023.acl-long.191. URL <https://aclanthology.org/2023.acl-long.191>.
- C.-Y. Lin. Rouge: A package for automatic evaluation of summaries. In *Text summarization branches out*, pages 74–81, 2004.
- S. Lin, B. Liu, J. Li, and X. Yang. Common diffusion noise schedules and sample steps are flawed. In *Proceedings of the IEEE/CVF winter conference on applications of computer vision*, pages 5404–5411, 2024.
- S. Liu, H. Lu, and J. Shao. Improved residual vector quantization for high-dimensional approximate nearest neighbor search. *arXiv preprint arXiv:1509.05195*, 2015.
- J. Lovelace, V. Kishore, Y. Chen, and K. Q. Weinberger. Diffusion guided language modeling. *ArXiv*, abs/2408.04220, 2024. URL <https://arxiv.org/pdf/2408.04220>.
- A. Lozhkov, L. Ben Allal, L. von Werra, and T. Wolf. Fineweb-edu, May 2024. URL <https://huggingface.co/datasets/HuggingFaceFW/fineweb-edu>.
- C. Lu, Y. Zhou, F. Bao, J. Chen, C. Li, and J. Zhu. Dpm-solver: A fast ode solver for diffusion probabilistic model sampling in around 10 steps. *Advances in Neural Information Processing Systems*, 35:5775–5787, 2022.
- A. Marfurt and J. Henderson. Sentence-level planning for especially abstractive summarization. In *ACL*, pages 1–14, 2021. URL <https://aclanthology.org/2021.newsum-1.1.pdf>.
- B. Minixhofer, J. Pfeiffer, and I. Vulić. Where’s the point? self-supervised multilingual punctuation-agnostic sentence segmentation. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 7215–7235, Toronto, Canada, July 2023. Association for Computational Linguistics. URL <https://aclanthology.org/2023.acl-long.398>.
- I. Mohr, M. Krimmel, S. Sturua, M. K. Akram, A. Koukounas, M. Günther, G. Mastrapas, V. Ravishankar, J. F. Martínez, F. Wang, Q. Liu, Z. Yu, J. Fu, S. Ognawala, S. Guzman, B. Wang, M. Werk, N. Wang, and H. Xiao. Multi-task contrastive learning for 8192-token bilingual text embeddings, 2024. URL <https://arxiv.org/abs/2402.17016>.
- N. Mostafazadeh, N. Chambers, X. He, D. Parikh, D. Batra, L. Vanderwende, P. Kohli, and J. Allen. A corpus and cloze evaluation for deeper understanding of commonsense stories. In *NAACL*, pages 839–849, June 2016. URL <https://aclanthology.org/N16-1098>.
- S. Narayan, S. B. Cohen, and M. Lapata. Don’t give me the details, just the summary! topic-aware convolutional neural networks for extreme summarization. *arXiv preprint arXiv:1808.08745*, 2018.
- J. Ni, G. H. Abrego, N. Constant, J. Ma, K. B. Hall, D. Cer, and Y. Yang. Sentence-t5: Scalable sentence encoders from pre-trained text-to-text models. *arXiv preprint arXiv:2108.08877*, 2021.
- A. Q. Nichol and P. Dhariwal. Improved denoising diffusion probabilistic models. In *International conference on machine learning*, pages 8162–8171. PMLR, 2021.
- A. Q. Nichol, P. Dhariwal, A. Ramesh, P. Shyam, P. Mishkin, B. McGrew, I. Sutskever, and M. Chen. GLIDE: Towards photorealistic image generation and editing with text-guided diffusion models. In K. Chaudhuri, S. Jegelka, L. Song, C. Szepesvari, G. Niu, and S. Sabato, editors, *Proceedings of the 39th International Conference on Machine Learning*, volume 162 of *Proceedings of Machine Learning Research*, pages 16784–16804. PMLR, 17–23 Jul 2022. URL <https://proceedings.mlr.press/v162/nichol22a.html>.

- M. Ning, M. Li, J. Su, A. A. Salah, and I. O. Ertugrul. Elucidating the exposure bias in diffusion models. *arXiv preprint arXiv:2308.15321*, 2023.
- NLLB Team, M. R. Costa-jussà, J. Cross, O. Çelebi, M. Elbayad, K. Heafield, K. Heffernan, E. Kalbassi, J. Lam, D. Licht, J. Maillard, A. Sun, S. Wang, G. Wenzek, A. Youngblood, B. Akula, L. Barrault, G. Mejjia-Gonzalez, P. Hansanti, J. Hoffman, S. Jarrett, K. R. Sadagopan, D. Rowe, S. Spruit, C. Tran, P. Andrews, N. F. Ayan, S. Bhosale, S. Edunov, A. Fan, C. Gao, V. Goswami, F. Guzmán, P. Koehn, A. Mourachko, C. Ropers, S. Saleem, H. Schwenk, and J. Wang. No language left behind: Scaling human-centered machine translation, 2022. URL <https://arxiv.org/abs/2207.04672>.
- OpenAI. GPT-4 technical report. *ArXiv*, abs/2303.08774, 2024. URL <https://arxiv.org/pdf/2303.08774>.
- K. Papineni, S. Roukos, T. Ward, and W.-J. Zhu. Bleu: a method for automatic evaluation of machine translation. In *Proceedings of the 40th Annual Meeting of the Association for Computational Linguistics*, pages 311–318, Philadelphia, Pennsylvania, USA, July 2002. Association for Computational Linguistics. doi: 10.3115/1073083.1073135. URL <https://aclanthology.org/P02-1040>.
- A. Parola, J. M. Lin, A. Simonsen, V. Bliksted, Y. Zhou, H. Wang, L. Inoue, K. Koelkebeck, and R. Fusaroli. Speech disturbances in schizophrenia: Assessing cross-linguistic generalizability of nlp automated measures of coherence. *Schizophrenia Research*, 259:59–70, 2023.
- W. Peebles and S. Xie. Scalable diffusion models with transformers. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 4195–4205, 2023.
- E. Perez, F. Strub, H. De Vries, V. Dumoulin, and A. Courville. Film: Visual reasoning with a general conditioning layer. In *Proceedings of the AAAI conference on artificial intelligence*, 2018.
- A. Radford, J. Wu, R. Child, D. Luan, D. Amodei, I. Sutskever, et al. Language models are unsupervised multitask learners. *OpenAI blog*, 1(8):9, 2019.
- C. Raffel, N. Shazeer, A. Roberts, K. Lee, S. Narang, M. Matena, Y. Zhou, W. Li, and P. J. Liu. Exploring the limits of transfer learning with a unified text-to-text transformer. *arXiv e-prints*, 2019.
- C. Raffel, N. Shazeer, A. Roberts, K. Lee, S. Narang, M. Matena, Y. Zhou, W. Li, and P. J. Liu. Exploring the limits of transfer learning with a unified text-to-text transformer. *Journal of Machine Learning Research*, 21(140):1–67, 2020. URL <http://jmlr.org/papers/v21/20-074.html>.
- N. Reimers and I. Gurevych. Sentence-bert: Sentence embeddings using siamese Bert-networks. *arXiv preprint arXiv:1908.10084*, 2019.
- R. Rombach, A. Blattmann, D. Lorenz, P. Esser, and B. Ommer. High-resolution image synthesis with latent diffusion models. *CoRR*, abs/2112.10752, 2021. URL <https://arxiv.org/abs/2112.10752>.
- P. K. Rubenstein, C. Asawaroengchai, D. D. Nguyen, A. Bapna, Z. Borsos, F. de Chaumont Quitry, P. Chen, D. E. Badawy, W. Han, E. Kharitonov, H. Muckenhirn, D. Padfield, J. Qin, D. Rozenberg, T. N. Sainath, J. Schalkwyk, M. Sharifi, M. T. Ramanovich, M. Tagliasacchi, A. Tudor, M. Velimirovic, D. Vincent, J. Yu, Y. Wang, V. Zayats, N. Zeghidour, Y. Zhang, Z. Zhang, L. Zilka, and C. H. Frank. Audiopalm: A large language model that can speak and listen. *CoRR*, abs/2306.12925, 2023. URL <https://doi.org/10.48550/arXiv.2306.12925>.
- T. Salimans and J. Ho. Progressive distillation for fast sampling of diffusion models. *arXiv preprint arXiv:2202.00512*, 2022.
- Seamless Communication, L. Barrault, Y.-A. Chung, M. C. Meglioli, D. Dale, N. Dong, M. Duppenthaler, P.-A. Duquenne, B. Ellis, H. Elsahar, J. Haaheim, J. Hoffman, M.-J. Hwang, H. Inaguma, C. Klaiiber, I. Kulikov, P. Li, D. Licht, J. Maillard, R. Mavlyutov, A. Rakotoarison, K. R. Sadagopan, A. Ramakrishnan, T. Tran, G. Wenzek, Y. Yang, E. Ye, I. Evtimov, P. Fernandez, C. Gao, P. Hansanti, E. Kalbassi, A. Kallet, A. Kozhevnikov, G. M. Gonzalez, R. S. Roman, C. Touret, C. Wong, C. Wood, B. Yu, P. Andrews, C. Balioglu, P.-J. Chen, M. R. Costa-jussà, M. Elbayad, H. Gong, F. Guzmán, K. Heffernan, S. Jain, J. Kao, A. Lee, X. Ma, A. Mourachko, B. Peloquin, J. Pino, S. Popuri, C. Ropers, S. Saleem, H. Schwenk, A. Sun,

- P. Tomasello, C. Wang, J. Wang, S. Wang, and M. Williamson. Seamless: Multilingual expressive and streaming speech translation. *ArXiv*, abs/2312.05187, 2023a. URL <https://arxiv.org/abs/2312.05187>.
- Seamless Communication, L. Barrault, Y.-A. Chung, M. C. Meglioli, D. Dale, N. Dong, P.-A. Duquenne, H. Elsahar, H. Gong, K. Heffernan, J. Hoffman, C. Klaiber, P. Li, D. Licht, J. Maillard, A. Rakotoarison, K. R. Sadagopan, G. Wenzek, E. Ye, B. Akula, P.-J. Chen, N. E. Hachem, B. Ellis, G. M. Gonzalez, J. Haaheim, P. Hansanti, R. Howes, B. Huang, M.-J. Hwang, H. Inaguma, S. Jain, E. Kalbassi, A. Kallet, I. Kulikov, J. Lam, D. Li, X. Ma, R. Mavlyutov, B. Peloquin, M. Ramadan, A. Ramakrishnan, A. Sun, K. Tran, T. Tran, I. Tufanov, V. Vogeti, C. Wood, Y. Yang, B. Yu, P. Andrews, C. Balioglu, M. R. Costa-jussà, O. Celebi, M. Elbayad, C. Gao, F. Guzmán, J. Kao, A. Lee, A. Mourachko, J. Pino, S. Popuri, C. Ropers, S. Saleem, H. Schwenk, P. Tomasello, C. Wang, J. Wang, and S. Wang. SeamlessM4T - massively multilingual & multimodal machine translation, 2023b. URL <https://arxiv.org/abs/2308.11596>.
- A. Shabalin, V. Meshchaninov, E. Chimbulatov, V. Lapikov, R. Kim, G. Bartosh, D. Molchanov, S. Markov, and D. Vetrov. Tencdm: Understanding the properties of diffusion model in the space of language model encodings. *ArXiv*, abs/2402.19097, 2024. URL <https://arxiv.org/pdf/2402.19097>.
- N. Shazeer. Glu variants improve transformer. *arXiv preprint arXiv:2002.05202*, 2020.
- J. Song, C. Meng, and S. Ermon. Denoising diffusion implicit models. *CoRR*, abs/2010.02502, 2020. URL <https://arxiv.org/abs/2010.02502>.
- N. Srivastava, D. Kuchelev, T. M. Ngoli, K. Shetty, M. Röder, D. Moussallem, H. Zahera, and A.-C. N. Ngomo. Lola—an open-source massively multilingual large language model. *arXiv preprint arXiv:2409.11272*, 2024.
- S. Sturua, I. Mohr, M. K. Akram, M. Günther, B. Wang, M. Krimmel, F. Wang, G. Mastrapas, A. Koukounas, N. Wang, and H. Xiao. jina-embeddings-v3: Multilingual embeddings with task lora, 2024. URL <https://arxiv.org/abs/2409.10173>.
- H. Su, W. Shi, J. Kasai, Y. Wang, Y. Hu, M. Ostendorf, W.-t. Yih, N. A. Smith, L. Zettlemoyer, and T. Yu. One embedder, any task: Instruction-finetuned text embeddings. *arXiv preprint arXiv:2212.09741*, 2022.
- J. Su, M. Ahmed, Y. Lu, S. Pan, W. Bo, and Y. Liu. Roformer: Enhanced transformer with rotary position embedding. *Neurocomputing*, 568:127063, 2024.
- X. Sun, Z. Sun, Y. Meng, J. Li, and C. Fan. Summarize, outline, and elaborate: Long-text generation via hierarchical supervision from extractive summaries. In *Proceedings of the 29th International Conference on Computational Linguistics*, pages 6392–6402, Gyeongju, Republic of Korea, Oct. 2022. International Committee on Computational Linguistics. URL <https://aclanthology.org/2022.coling-1.556>.
- Team GLM, A. Zeng, B. Xu, B. Wang, C. Zhang, D. Yin, D. Zhang, D. Rojas, G. Feng, H. Zhao, H. Lai, H. Yu, H. Wang, J. Sun, J. Zhang, J. Cheng, J. Gui, J. Tang, J. Zhang, J. Sun, J. Li, L. Zhao, L. Wu, L. Zhong, M. Liu, M. Huang, P. Zhang, Q. Zheng, R. Lu, S. Duan, S. Zhang, S. Cao, S. Yang, W. L. Tam, W. Zhao, X. Liu, X. Xia, X. Zhang, X. Gu, X. Lv, X. Liu, X. Liu, X. Yang, X. Song, X. Zhang, Y. An, Y. Xu, Y. Niu, Y. Yang, Y. Li, Y. Bai, Y. Dong, Z. Qi, Z. Wang, Z. Yang, Z. Du, Z. Hou, and Z. Wang. ChatGLM: A family of large language models from glm-130b to glm-4 all tools. *ArXiv*, abs/2406.12793, 2024. URL <https://arxiv.org/pdf/2406.12793>.
- The Llama3 team. The Llama 3 herd of models. *ArXiv*, abs/2407.21783, 2024. URL <https://arxiv.org/pdf/2407.21783>.
- H. Touvron, L. Martin, K. Stone, P. Albert, A. Almahairi, Y. Babaei, N. Bashlykov, S. Batra, P. Bhargava, S. Bhosale, D. Bikel, L. Blecher, C. C. Ferrer, M. Chen, G. Cucurull, D. Esiobu, J. Fernandes, J. Fu, W. Fu, B. Fuller, C. Gao, V. Goswami, N. Goyal, A. Hartshorn, S. Hosseini, R. Hou, H. Inan, M. Kardas, V. Kerkez, M. Khabsa, I. Kloumann, A. Korenev, P. S. Koura, M.-A. Lachaux, T. Lavril, J. Lee, D. Liskovich, Y. Lu, Y. Mao, X. Martinet, T. Mihaylov, P. Mishra, I. Molybog, Y. Nie, A. Poulton, J. Reizenstein, R. Rungta, K. Saladi, A. Schelten, R. Silva, E. M. Smith, R. Subramanian, X. E. Tan, B. Tang, R. Taylor, A. Williams, J. X. Kuan, P. Xu, Z. Yan, I. Zarov, Y. Zhang, A. Fan, M. Kambadur, S. Narang, A. Rodriguez, R. Stojnic, S. Edunov, and T. Scialom. Llama 2: Open foundation and fine-tuned chat models, 2023.

- T. Ubukata, J. Li, and K. Tei. Diffusion model for planning: A systematic literature review. *ArXiv*, abs/2408.10266, 2024. URL <https://arxiv.org/pdf/2408.10266>.
- A. Üstün, V. Aryabumi, Z.-X. Yong, W.-Y. Ko, D. D’souza, G. Onilude, N. Bhandari, S. Singh, H.-L. Ooi, A. Kayid, et al. Aya model: An instruction finetuned open-access multilingual language model. *arXiv preprint arXiv:2402.07827*, 2024.
- C. Wang, S. Chen, Y. Wu, Z. Zhang, L. Zhou, S. Liu, Z. Chen, Y. Liu, H. Wang, J. Li, L. He, S. Zhao, and F. Wei. Neural codec language models are zero-shot text to speech synthesizers. *CoRR*, abs/2301.02111, 2023. URL <https://doi.org/10.48550/arXiv.2301.02111>.
- L. Wang, N. Yang, X. Huang, L. Yang, R. Majumder, and F. Wei. Improving text embeddings with large language models. In *ACL*, pages 11897–11916, 2024a. URL <https://aclanthology.org/2024.acl-long.642>.
- L. Wang, N. Yang, X. Huang, L. Yang, R. Majumder, and F. Wei. Multilingual e5 text embeddings: A technical report, 2024b. URL <https://arxiv.org/abs/2402.05672>.
- A. Warstadt, A. Singh, and S. R. Bowman. Neural network acceptability judgments. *Transactions of the Association for Computational Linguistics*, 7:625–641, 2019. doi: 10.1162/tacl_a_00290. URL <https://aclanthology.org/Q19-1040>.
- S. Welleck, I. Kulikov, S. Roller, E. Dinan, K. Cho, and J. Weston. Neural text generation with unlikelihood training. *arXiv preprint arXiv:1908.04319*, 2019.
- Y. Yang, G. H. Abrego, S. Yuan, M. Guo, Q. Shen, D. Cer, Y.-H. Sung, B. Strope, and R. Kurzweil. Improving multilingual sentence embedding using bi-directional dual encoder with additive margin softmax. *arXiv preprint arXiv:1902.08564*, 2019.
- J. Ye, J. Gao, S. Gong, L. Zheng, X. Jiang, Z. Li, and L. Kong. Beyond autoregression: Discrete diffusion for complex reasoning and planning. *ArXiv*, abs/2410.14157, 2024. URL <https://arxiv.org/pdf/2410.14157>.
- Y. Yin, J. Ding, K. Song, and Y. Zhang. Semformer: Transformer language models with semantic planning. *ArXiv*, abs/2409.11143, 2024. URL <https://arxiv.org/pdf/2409.11143>.
- N. Zeghidour, A. Luebs, A. Omran, J. Skoglund, and M. Tagliasacchi. Soundstream: An end-to-end neural audio codec. *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, 30:495–507, 2021.
- B. Zhang and R. Sennrich. Root mean square layer normalization. *Advances in Neural Information Processing Systems*, 32, 2019.
- X. Zhang, Y. Zhang, D. Long, W. Xie, Z. Dai, J. Tang, H. Lin, B. Yang, P. Xie, F. Huang, M. Zhang, W. Li, and M. Zhang. mgte: Generalized long-context text representation and reranking models for multilingual text retrieval, 2024. URL <https://arxiv.org/abs/2407.19669>.
- Y. Zhang, J. Gu, Z. Wu, S. Zhai, J. Susskind, and N. Jaitly. Planner: Generating diversified paragraph via latent language diffusion model. In *NeurIPS*, 2023.

A Technical consideration for data preparation

Since our modeling approach uses a fixed encoder and a fixed document segmentation method, we decided to use pre-computed SONAR embeddings instead of producing them on-the-fly for each training run. This allows for faster iteration on the same data mix, trading expensive GPU compute against storage capacity.

As we are storing sequences of SONAR embedding, which are fixed size tensors of 1024 floats, the storage requirements become more demanding than storing the raw text. For one terra bytes of raw text data we need to store between fifteen and twenty terra bytes of encoded data. Overall, this trade-off in space vs compute reduces the GPU memory occupation and the compute load and lets use iterate faster. Typically, with on single GPU we can produce around 300-400 SONAR sentence embeddings per second whereas by loading precomputed data (potentially from remote storage) we can load over 20 thousand embeddings per second per GPU (with around 15 CPU per GPU).

We store sequences of SONAR embeddings with 16 bits precision (FP16) in parquet datasets. Embeddings remain aligned with the segmented texts and the parquet binary format and library ecosystem is well suited for storing and loading efficiently such complex data structures. Parquet also lets us store extra data (such as quality metrics for each sentences) and enables non-trivial last mile data filtering and transformation.

For training the LCM, we processed around four billion documents, generating 310 billion sentences with an average of 27 tokens per sentences for 88 characters length on average; totaling a bit more than 889 terra-bytes of raw text.

B Open Sourced Code

In the spirit of reproducibility, we release under an open source license the training, evaluation and data processing code for the LCM. This is available at https://github.com/facebookresearch/large_concept_model.

The training code is based on the FAIRSEQ2 framework (Balioglu, 2023) that allowed us to build and iterate over the different model architectures discussed above. While FAIRSEQ2 shares the same name as the popular fairseq toolchain, its API architecture is different. It is not a monolithic toolchain but a set of modules that can be composed, this allows us to build different architectures side by side and easily share training components.

We also release our evaluation framework so the evaluation tasks reported in 3.1 and comparison between the LCM and other models can easily be reproduced. The evaluation framework provides a clear abstraction between *predictors*, *tasks* and *data loading*, which again, makes it modular and lets us describe a set of tasks to be evaluated. The evaluation framework can be run locally or distributed over a SLURM cluster to run evaluations at scale.

Finally, we release an updated version of stopes¹⁶ to simplify large scale data pre-processing on a SLURM cluster. This was used to run the sentence segmentation and SONAR encoding described in section 2.2. The stopes data processing framework deals with scheduling and monitoring large number of jobs on SLURM or to run everything locally for small scale jobs. It provides an API compatible with `ray.data`¹⁷ that makes it easy to process large datasets in blocks and apply transform

¹⁶<https://github.com/facebookresearch/stopes>

¹⁷<https://docs.ray.io/>

function over it. This makes our code reusable outside of a SLURM cluster as it can also be used with a ray.io cluster.

C System prompt: Generation of Topic Descriptions

You are a topic description generator. Your job is to read an extract of text and then generate a topic description. The extract may be well formed or not. The topic description you will write will be at most one sentence in length, and use as few words as possible. However, it can not be generic and it can not contain any profanity.

Here is an example of an extract, an ideal topic description, and some examples of bad topic descriptions:

Example extract: "One day, one neighborhood of the city was completely devastated. Glass windows were shattered, shops turned upside down, and many civilian killed. Superman instantly recognized the signature of one of his old enemies, Voltar, who he had barely beaten in the past. This was a message to him: "I challenge you! Come find me!""

An example of a good topic description: An old enemy of Superman's, Voltar, appeared and challenged him.

An example of a bad topic description: Superman

An example of a bad topic description: Voltar

An example of a bad topic description:

D User prompt: LLM As a Judge - Coherence

Below is a text extract. Your task is to analyze the extract and assign a coherence score between 0 and 5 inclusive, where:

0: The text is completely incoherent and lacks any logical connection.

1: The text has some minor connections, but overall it is disjointed and hard to follow.

2: The text has some coherence, but it is still difficult to understand due to unclear relationships between ideas.

3: The text is moderately coherent, with some clear connections between ideas, but may lack depth or clarity.

4: The text is highly coherent, with clear and logical connections between ideas, making it easy to follow.

5: The text is extremely coherent, with a clear and concise structure, making it effortless to understand.

You will provide a score ONLY. Do NOT also provide an explanation.

The extract: <extract>

After examining the extract, the coherence score between 0 and 5 inclusive is: