

## Title

Wearable intelligent throat enables natural speech in stroke patients with dysarthria

## Authors

Chenyu Tang<sup>†1</sup>, Shuo Gao<sup>†\*2</sup>, Cong Li<sup>†2</sup>, Wentian Yi<sup>1</sup>, Yuxuan Jin<sup>3</sup>, Xiaoxue Zhai<sup>4</sup>, Sixuan Lei<sup>5</sup>, Hongbei Meng<sup>2</sup>, Zibo Zhang<sup>1</sup>, Muzi Xu<sup>1</sup>, Shengbo Wang<sup>2</sup>, Xuhang Chen<sup>6</sup>, Chenxi Wang<sup>2</sup>, Hongyun Yang<sup>2</sup>, Ningli Wang<sup>7</sup>, Wenyu Wang<sup>8</sup>, Jin Cao<sup>9</sup>, Xiaodong Feng<sup>10</sup>, Peter Smielewski<sup>6</sup>, Yu Pan<sup>4</sup>, Wenhui Song<sup>11</sup>, Martin Birchall<sup>12</sup>, and Luigi G. Occhipinti<sup>\*1</sup>

## Affiliations

<sup>1</sup>Department of Engineering, University of Cambridge, Cambridge, UK

<sup>2</sup>School of Instrumentation and Optoelectronic Engineering, Beihang University, Beijing, China

<sup>3</sup>Cavendish Laboratory, University of Cambridge, Cambridge, UK

<sup>4</sup>Department of Rehabilitation Medicine, Beijing Tsinghua Changgung Hospital, Tsinghua University, Beijing, China

<sup>5</sup>Shenzhen International Graduate School, Tsinghua University, Shenzhen, China

<sup>6</sup>Department of Clinical Neurosciences, University of Cambridge, Cambridge, UK

<sup>7</sup>Beijing Tongren Hospital, Capital Medical University, Beijing, China

<sup>8</sup>Thrust of Smart Manufacturing, Hong Kong University of Science and Technology (Guangzhou)

<sup>9</sup>School of Life Sciences, Beijing University of Chinese Medicine, Beijing, China

<sup>10</sup>Department of Rehabilitation Center, The First Affiliated Hospital of Henan University of Chinese Medicine, Zhengzhou, China

<sup>11</sup>Department of Surgical Biotechnology, University College London, London, United Kingdom

<sup>12</sup>Royal National Ear Nose and Throat and Eastman Dental Hospitals, University College London Hospital, London, UK

<sup>†</sup>These authors contributed equally: Chenyu Tang, Shuo Gao, Cong Li

<sup>\*</sup>Correspondence to: Shuo Gao (shuo\_gao@buaa.edu.cn) and Luigi G. Occhipinti (lgo23@cam.ac.uk)

## Abstract

Wearable silent speech systems hold significant potential for restoring communication in patients with speech impairments. However, seamless, coherent speech remains elusive, and clinical efficacy is still unproven. Here, we present an AI-driven intelligent throat (IT) system that integrates throat muscle vibrations and carotid pulse signal sensors with large language model (LLM) processing to enable fluent, emotionally expressive communication. The system utilizes ultrasensitive textile strain sensors to capture high-quality signals from the neck area and supports token-level processing for real-time, continuous speech decoding, enabling seamless, delay-free communication. In tests with five stroke patients with dysarthria, IT's LLM agents intelligently corrected token errors and enriched sentence-level emotional

and logical coherence, achieving low error rates (4.2% word error rate, 2.9% sentence error rate) and a 55% increase in user satisfaction. This work establishes a portable, intuitive communication platform for patients with dysarthria with the potential to be applied broadly across different neurological conditions and in multi-language support systems.

## I. Main

Neurological diseases such as stroke, amyotrophic lateral sclerosis (ALS), and Parkinson’s disease frequently result in dysarthria—a severe motor-speech disorder that compromises neuromuscular control over the vocal tract. This impairment drastically restricts effective communication, lowers quality of life, substantially impedes the rehabilitation process, and can even lead to severe psychological issues [1, 2, 3, 4]. Augmentative and alternative communication (AAC) technologies have been developed to address these challenges, including letter-by-letter spelling systems utilizing head or eye tracking [5, 6, 7, 8] and neuroprosthetics powered by brain-computer interface (BCI) devices [9, 10, 11, 12]. While head or eye tracking systems are relatively straightforward to implement, they suffer from slow communication speeds. Neuroprosthetics, while transformative for severe paralysis cases, often rely on invasive, complex recordings and processing of neural signals. For individuals retaining partial control over laryngeal or facial muscles, a strong need remains for solutions that are more intuitive and portable (SNote 1).

A promising solution lies in wearable silent speech devices that capture non-acoustic signals, such as subtle skin vibrations [13, 14, 15, 16, 17] or electrophysiological signals from the speech motor cortex [18, 19, 20, 21]. These technologies offer non-invasiveness, comfort, and portability, with potential for seamless daily integration. Yet, despite their promise, current systems remain in their infancy, achieving reliable, discrete word decoding in healthy users but showing limited success in patient trials [13, 14, 15]. More critically, these systems fall short of delivering truly natural communication—requiring both delay-free expression and consistent contextual coherence, capabilities essential for fully effective and meaningful interactions.

To advance wearable silent speech systems for real-world dysarthria patient use, we developed an AI-driven intelligent throat (IT) system that captures extrinsic laryngeal muscle vibrations and carotid pulse signals, integrating silent speech and emotional states analysis in real-time. The system generates personalized, contextually appropriate sentences that accurately reflect patients' intended meaning (Figure 1). It employs ultrasensitive textile strain sensors, fabricated using advanced printing techniques, to ensure comfortable, durable, and high-quality signal acquisition [14, 22]. By analyzing speech signals at the token level (~100ms), our approach outperforms traditional time-window methods, enabling continuous, fluent word and sentence expression in real time. Knowledge distillation further reduces computational latency by 76%, significantly enhancing communication fluidity. Large language models (LLMs) serve as intelligent agents, automatically correcting token classification errors and generating personalized, context-aware speech by integrating emotional states and environmental cues. Pre-trained on a dataset from 10 healthy individuals, the system achieved a word error rate (WER) of 4.2% and a sentence error rate (SER) of 2.9% when fine-tuned on data from five dysarthric stroke patients. Additionally, the integration of emotional states and contextual cues further personalizes and enriches the decoded sentences, resulting in a 55% increase in user satisfaction and enabling dysarthria patients to communicate with fluency and naturalness comparable to that of healthy individuals. STable 1 provides a comprehensive comparison between the IT system and state-of-the-art wearable silent speech systems.

## II. Results

### The intelligent throat system

The IT system consists primarily of hardware (a smart choker embedding textile strain sensors and a wireless readout printed circuit board (PCB)) and software components (machine learning models and LLM agents). Silent speech signals generated in real time by the user's silent expressions are decoded by a token decoding network and synthesized into an initial sentence by the token synthesis agent (TSA). Simultaneously, pulse signals are collected from the smart choker device and processed by an emotion decoding network to determine the user's real-time emotional status. The sentence expansion agent (SEA) intelligently expands the TSA-generated sentence, incorporating personalized emotion labels and objective contextual background data to produce a refined, emotionally expressive, and logically coherent sentence that captures the user's intended meaning (Fig. 1, SVideo 2). Each component of the IT system is elaborated upon in the following sections.

Fig. 2a shows the structure of the strain sensing choker screen-printed on an elastic knitted textile. The choker features two channels located at the front and side of the neck, designed to monitor the strain applied to the skin by the muscles near the throat and the carotid artery (SFig. 1). The graphene layer printed on the textile forms ordered cracks along the stress concentration areas of the textile lattice to detect subtle skin vibrations [14]. Silver electrodes are connected to the integrated PCB on the choker. A rigid strain isolation layer with high Young's modulus is printed around each channel to reduce crosstalk between the two channels and the variable strains caused by wearing. Due to the difference in Young's modulus between the elastic textile substrate and the strain isolation layer, less than 1% of external strain is transmitted to the interior when wearing the choker, while the internal sensing areas remain soft and elastic (SFig. 2) [22]. For uniaxial stretching from 1-10 Hz, the printed textile-based graphene strain sensor shows good linear behaviour, producing a response over 10% to subtle strains of 0.1%, and maintains a gauge factor (GF) over 100 during high-frequency stretching (Fig. 2b). Furthermore, our previous studies have confirmed the reliability of the printed textile-based strain sensors with high robustness, durability and washability, as well as high levels of comfort, biocompatibility and breathability [14, 22].

To operate the system and enable wireless communication between the IT choker and server, the PCB was designed for bi-channel measurements (i.e., silent speech and carotid pulse signals), enabling simultaneous acquisition of speech and emotional cues. The PCB integrates a low-power Bluetooth module (Fig. 2c) for continuous data transmission while optimizing energy efficiency for extended use. Key components of the PCB include an analog-to-digital converter (ADC) for high-fidelity signal digitization and a microcontroller unit (MCU) that manages data processing and transmission (Fig. 2e, SFig. 4, and SFig. 5). Power supply, operational amplifiers, and the reference voltage chip are configured to ensure stable signal amplification, catering to the sensitivity requirements of both strain and pulse sensors. For the energy management system, a comprehensive power budget analysis reveals that the designed PCB operates with a total power consumption of 76.5 mW (Fig. 2f). The main power-consuming components are the Bluetooth module (29.7 mW) and amplification circuits (31.9 mW). To extend operational time and support portable use, a 1800 mWh battery was incorporated, providing sufficient capacity for continuous operation throughout an entire day without recharging.

### Token-level speech decoding

Current wearable silent speech systems operate by recognizing discrete words or predefined sentences and lack the ability for continuous, real-time expression analysis typical of the human brain [45]. This limitation arises because these systems rely on fixed time windows (typically 1–3 seconds) for word decoding, requiring users to complete each word within a set interval and pause until the next window to continue [13-21]. Such constraints lead to fragmented expression and unnatural user experience. To address this, we developed a high-resolution tokenization method for signal segmentation (Fig. 2f), dividing speech signals into fine-grained ~100ms segments for continuous word label recognition. This granular segmentation ensures that each token accurately corresponds to a specific part of a single word and is labeled accordingly. This setup enables users to speak fluidly without worrying about timing constraints, as the system continuously classifies and aggregates tokens into coherent words and sentences. Our optimization determined that a token length of 144 ms offers the ideal balance: it minimizes boundary confusion from longer tokens while avoiding the increased computational demands associated with shorter tokens.

While high-resolution tokenization improves fluidity, shorter tokens inherently contain limited context, making them less effective for accurate word decoding. Temporal machine learning models, like recurrent neural networks (RNN) or transformers, could capture contextual dependencies, but their complexity and computational cost render them suboptimal for wearable silent speech systems [23, 24, 25], which prioritize real-time operation. To balance context awareness and computational efficiency, we implemented an explicit context augmentation strategy (Fig. 3a), where each sample consists of  $N$  tokens:  $N-1$  preceding tokens provide context, and the current token determines the sample's label. For initial tokens, any missing preceding tokens are padded with blank tokens to ensure completeness. We found  $N=15$  tokens to be optimal (Fig. 3c), with accuracy initially increasing as tokens accumulate, then declining due to insufficient context at lower counts and gradient decay or information loss at higher counts [26]. This strategy enables the use of efficient one-dimensional convolutional neural networks (1D-CNNs) instead of computationally intensive temporal models for token decoding [27, 28]. Attention maps reveal that signals from preceding regions indeed contribute to token decoding, validating the effectiveness of the explicit context augmentation strategy (SFig.10).

To further enhance model efficiency and accuracy on patients' data, we designed the training pipeline shown in Fig. 3b. The model was pre-trained on a larger dataset from healthy individuals and then fine-tuned on the limited patients' data, leveraging shared signal features to enhance patient-specific decoding. After only 25 repetitions per word in few-shot learning, the model achieved a token classification accuracy of 92.2% (Fig. 3d). In contrast, a model trained from scratch using solely patients' data could only reach an accuracy of 79.8%. Additionally, we employed response-based knowledge distillation [29] to transfer knowledge from a larger 1D ResNet-101 model to a smaller 1D ResNet-18, reducing computational load by 75.6% while maintaining high accuracy, with only a 0.9% drop from the teacher model, achieving 91.3% (Fig. 3e). Fig. 3f and Fig. 3g display the confusion matrix and UMAP feature visualization for token decoding [30]. Over 90% of the classification errors involved confusion between class 0 (blank tokens) and neighbouring word tokens. As shown in later analyses of the LLM agent's performance, such boundary errors can be effectively corrected during token-to-word synthesis by the token synthesis agent (TSA).

### **Decoding of emotional states**

To enrich sentence coherence by providing emotional context, we decode emotional states from carotid

pulse signals. Emotional state recognition can typically be achieved through a variety of methods, including analysis of facial images from cameras, audio speech signals, and various physiological indicators such as heart rate and blood pressure [31, 32, 33]. In line with our objective of creating a highly integrated wearable system, we chose carotid pulse signals as a biomarker for emotional decoding. Using 5-second windows, we segmented patients' pulse signals into samples to construct a dataset, focusing on three common emotion categories for stroke patients: neutral, relieved, and frustrated (data collection protocol detailed in Methods). Fig. 4a shows the discrete Fourier transform (DFT) distributions for each emotion, highlighting distinct frequency characteristics among these emotional states. Accordingly, we incorporated DFT frequency extraction into the decoding pipeline shown in Fig. 4b, where removal of the DC component, Z-score normalization, and DFT are sequentially applied before feeding the values into a classifier for categorization. Fig. 4c illustrates the performance of different classifiers with and without DFT frequency extraction. The results show a significant improvement in decoding accuracy with DFT. The optimal model was the 1D-CNN with DFT, achieving an accuracy of 83.2%, with its confusion matrix displayed in Fig. 4d. The SHAP values reveal that the emotion decoding model primarily focuses on low-frequency signals in the 0-2 Hz range, which is consistent with the pulse signal range demonstrated by the DFT (SFig. 11).

In addition to the silent speech and carotid pulse signals analyzed in this study, various physiological activities generate distinct vibrational signals in the neck area, which can introduce artefacts hindering analysis [34, 35]. Fig. 4e shows the frequency and magnitude distributions of several prominent signals in this region. Our observations revealed that silent speech exhibits a relatively strong magnitude, and in applications with the IT, vibration can propagate transversely from the throat center to the carotid artery, introducing crosstalk in the pulse signal. Due to the considerable frequency overlap between silent speech and pulse signals, digital filters are non-ideal for effective artefacts suppression [36]. While adding reference channels could theoretically help, it does not align with the goal of a highly integrated IT [37]. To address this issue, we employed a stress isolation treatment using a polyurethane acrylate (PUA) layer, as shown in Fig. 2a, to prevent strain crosstalk propagation along the IT. The theoretical basis of this isolation strategy has been thoroughly discussed in our previous study [22]. Fig. 4f compares pulse signals with and without strain isolation treatment when silent speech occurs concurrently (the vowel "a" introduced at 2.5s), demonstrating significant crosstalk resilience in the treated IT.

## **LLM agents for sentence synthesis and intelligent expansion**

To naturally and coherently synthesize sentences that accurately reflect the patient's intended expression from the decoded token and emotion labels, we introduced two LLM agents based on the GPT-4o-mini API (Fig. 5a): the token synthesis agent (TSA) and the sentence expansion agent (SEA). The TSA merges token labels directly into words silently expressed by the patient and combines them into sentences (left). The SEA, on the other hand, leverages emotion labels and objective information, such as time and weather, to expand these basic sentences into logically coherent, personalized expressions that better capture the patient's true intent. Through a simple interaction (in this study, two consecutive nods), the IT system enables seamless switching between the direct output and the enriched, expanded sentence.

To optimize the performance of the TSA, we refined the prompt design [38]. First, we optimized the prompt length (Fig. 5b), observing a trend where both WER and SER improved with increasing prompt length up to 400 words before eventually deteriorating for higher lengths. We attribute this trend to the fact that longer prompts provide clearer synthesis instructions, but overly lengthy prompts dilute the

model's focus ability. Additionally, we compared performance with and without example cases, where the agent was provided with five examples of token label sequences and their corrected word outputs. Including examples significantly improved synthesis accuracy (Fig. 5c). Finally, we evaluated the effect of providing empirical constraints, which specify typical token counts for words of various lengths. Performance improved considerably when constraints were included. Under optimal prompt conditions, TSA achieved its best performance with a WER of 4.2% and an SER of 2.9%.

We also assessed and refined the performance of the SEA. Patient satisfaction with the expanded sentences was evaluated through a questionnaire (see STable 4 for criteria details). Following Chain-of-Thought (CoT) optimization [39] and the inclusion of patient-provided expansion examples, the expanded sentences scored significantly higher across multiple criteria (Fig. 5f). Contribution analysis revealed that emotion labels made a substantial impact on emotion accuracy, while objective information notably improved fluency, jointly contributing to the overall satisfaction with the expanded sentences compared to the basic word-only output (Fig. 5e). Under optimal prompt conditions, the SEA-generated expanded sentences resulted in a 55% increase in overall patient satisfaction compared to the TSA's direct output, raising satisfaction from "somewhat satisfied" to "fully satisfied" levels (SFig. 12 and SFig. 13).

In both operating modes, sentences generated by the TSA and SEA agents are sent to an open-source text-to-speech model [44], which synthesizes audio that matches the patient's natural voice for playback. In real-world applications, the delay between the completion of the user's silent expression and the sentence playback is approximately 1 second (SNote 2). This low latency effectively supports seamless and natural communication in practical settings.

### III. Discussion

In this work, we introduce the IT, an advanced wearable system designed to empower dysarthric stroke patients to communicate with the fluidity, intuitiveness, and expressiveness of natural speech. Comprehensive analysis and user feedback affirm the IT's high performance in fluency, accuracy, emotional expressiveness, and personalization. This success is rooted in its innovative design: ultrasensitive textile strain sensors capture rich and high-quality vibrational signals from the laryngeal muscles and carotid artery, while high-resolution tokenized segmentation enables users to communicate freely and continuously without expression delays. Additionally, the integration of LLM agents enables intelligent error correction and contextual adaptation, delivering exceptional decoding accuracy (WER < 5%, SER < 3%) and a 55% increase in user satisfaction. The IT thus sets a new benchmark in wearable silent speech systems, offering a naturalistic, user-centered communication aid.

Future efforts in several key areas will guide the continued development of the IT system. First, expanding its adaptability to a wider range of neurological conditions and demographic groups will make the technology more inclusive. Second, enhancing its linguistic diversity and multilingual support will allow for more personalized communication across language barriers. Finally, miniaturizing the system within an edge computing framework will facilitate seamless integration into real-world settings, boosting usability and accessibility.

Looking ahead, the advantages of the IT extend beyond enhancing everyday communication; they contribute to the holistic health of neurological patients, encompassing both physical and psychological well-being. The regained fluency in communication allows patients to re-engage in social interactions,

reducing isolation and the associated risk of depression. Moreover, effective communication facilitates real-time, personalized adjustments by rehabilitation therapists, supporting patients' recovery from motor impairments like hemiplegia. Together, these capabilities position the IT as a comprehensive tool for restoring independence and improving quality of life for individuals with neurological conditions.

## **IV. Methods**

### **Materials**

TIMREX KS 25 Graphite (particle size of 25  $\mu$  m) was sourced from IMERYS. Stretchable conductive silver ink was obtained from Dycotec Materials Ltd. Ethyl cellulose was purchased from SIGMA-ALDRICH. Flexible UV Resin Clear was acquired from Photocentric Ltd. The textile substrate, composed of 95% Polyester and 5% spandex, was procured from Jelly Fabrics Ltd.

### **Ink formulation**

The graphene ink for screen printing was prepared following a reported method. Briefly, 100g of graphite powder and 2g of ethyl cellulose (EC) were mixed in 1L of isopropyl alcohol (IPA) and stirred at 3000 rpm for 30 minutes. The mixture was then added into a high-pressure homogenizer (PSI-40) at 2000 bar pressure for 50 cycles to obtain graphene dispersion. The graphene dispersion is centrifuged at 5000g for 30 min to remove unexfoliated graphite.

### **Fabrication of textile strain sensor**

The textile substrate was washed with detergent, thoroughly dried, and then treated with UV-ozone for 5 minutes to clean the surface. Screen printing was performed using a 165T polyester silk screen on a semi-automatic printer (Kippax & Sons Ltd.) set with a squeegee angle of 45 degrees, a spacer of 2 mm, a coating speed of 10 mm/s, and a printing speed of 40 mm/s. Graphene ink, silver paste, and PUA were successively printed to form the sensing layer, electrodes, and strain isolation layer, respectively. After printing the PUA, the textile was exposed to UV light for 5 minutes. After each printing pass, the textile was air-dried. Following printing, the sensor was dried at 80 °C overnight. A biaxial strain of approximately 10% was then applied to induce the formation of ordered cracks.

### **Characterization**

Scanning Electron Microscopy (SEM) images were taken with a Magellan 400, after sputtering the textile samples with a 5 nm layer of gold to enhance conductivity. Optical images were captured using an Olympus microscope. Tensile properties of the textile strain sensors were evaluated using a Deben Microtest 200N Tensile Stage and an INSTRON universal testing system. Electrical signals were recorded concurrently with a potentiostat (EmStat4X, PalmSens) and a multiplexer (MUX, PalmSens). Copper tape was crimped onto the contact pads of the samples, supplemented with a small amount of silver paste to improve electrical contact.

### **Wireless PCB for data readout**

A custom wireless PCB was developed for efficient, continuous data acquisition and transmission within the IT system. Powered by a TP4065 lithium charger and a 3.3V regulator, the PCB ensures stable operation via battery or USB. The STM32G431 microcontroller captures silent speech and carotid pulse

signals through two ADC channels, with an OPA2192 operational amplifier for high-precision signal conditioning, amplifying low-level signals and enhancing overall data fidelity. A BLE module (BLE-SER-A-ANT) transmits real-time data via UART, enabling seamless, delay-free communication.

### **Silent speech data acquisition**

We recruited 10 healthy subjects (mean age:  $25.3 \pm 4.1$  years; 6 males, 4 females) and 5 stroke patients with dysarthria (mean age:  $43.9 \pm 8.3$  years; 4 males, 1 female) for silent speech signal collection, in compliance with Ethics Committee approval from the First Affiliated Hospital of Henan University of Chinese Medicine, approval no. 2023HL-142-01. A corpus was developed consisting of 47 Chinese words commonly used by stroke patients in daily communication, along with 20 sentences constructed from these words (see STable 2 and STable 3). For the healthy subject dataset, we collected 100 repetitions per word and 50 repetitions per sentence. For the patient dataset, we gathered 50 repetitions per word and 50 per sentence.

The healthy subject data serves as a critical baseline for initial model training, enabling the model to establish foundational patterns in silent speech signals. This pre-training facilitates improved generalization and performance when later fine-tuning the model on the limited data from dysarthric patients, ultimately enhancing decoding accuracy and robustness in patient-specific applications. The silent speech signals were segmented into tokens at 144 ms intervals. Each token was combined with the preceding 14 tokens to form a sample, allowing the model to incorporate context. The sample's label corresponds to the word of the current token. The signals were originally recorded at a sampling rate of 10 kHz and subsequently downsampled to 1 kHz before tokenization. Before neural network analysis, each sample was uniformly preprocessed with detrending and z-score normalization.

### **Protocol for emotion data collection**

Emotional pulse data was collected concurrently with silent speech signals, ensuring synchronized datasets that capture both speech-related and underlying physiological responses. To achieve accurate labeling, each emotion—neutral, relieved, and frustrated—was elicited through a carefully structured protocol involving audio-induced emotional states [40, 41, 42]. The emotions were induced via the international affective digitized sounds (2nd Edition; IADS-2) [43]. The three emotions were chosen as they are the most frequently encountered emotions in dysarthric patients' daily communication. Labeling was verified through collaboration between the participants and the therapist to ensure the successful and reliable induction of each target emotion. To balance sufficient information within each window and achieve the necessary resolution for emotion detection, pulse signals were segmented into 5-second samples. A 50% window overlap was applied to increase the training set size, enhancing model learning and generalization. The signals were originally recorded at a sampling rate of 10 kHz and subsequently downsampled to 200 Hz before analysis.

### **Software environment and model training**

Signal preprocessing was performed on a MacBook Pro equipped with an M1 Max CPU. Network training was conducted using Python 3.8.13, Miniconda 3, and PyTorch 2.0.1 in a performance-optimized environment. Training acceleration was enabled by CUDA on NVIDIA A100 GPU. The detailed training parameters for all models can be found in SFig. 8 and SFig. 9.



## Data availability

The datasets supporting this study will be available from the GitHub repository before publication.

## Code availability

The code supporting this study will be available from the GitHub repository before publication.

## Acknowledgments

This work was partially supported by the British Council (Grant Contract No. 45371261), the UK Engineering and Physical Science Research Council (EPSRC, grants No. EP/K03099X/1, EP/W024284/1) and Haleon through the CAPE partnership contract (University of Cambridge Ref. No. G110480).

## References

- [1] Enderby, P. Disorders of communication. *Neurological Rehabilitation* **110**, 273–281 (2013).
- [2] Tang, C. *et al.* A roadmap for the development of human body digital twins. *Nature Reviews Electrical Engineering* **1**, 199–207 (2024)
- [3] Zinn, S., *et al.* The effect of poststroke cognitive impairment on rehabilitation process and functional outcome. *Archives of physical medicine and rehabilitation* **85**, 1084–1090 (2004).
- [4] Teshaboeva, F. Literacy education of speech impaired children as a pedagogical psychological problem." *Confrencea* **5**, 299–302 (2023).
- [5] Ju, X. *et al.* A systematic review on voiceless patients' willingness to adopt high-technology augmentative and alternative communication in intensive care units. *Intensive and Critical Care Nursing* **63**, 102948 (2020).
- [6] Megalingam, R. *et al.* Sakthiprasad Kuttankulungara Manoharan, Gokul Riju & Sreekanth Makkal Mohandas. NETRAVAAD: Interactive Eye Based Communication System For People With Speech Issues. *IEEE Access* **12**, 69838–69852 (2024).
- [7] Ezzat, M. *et al.* Blink-To-Live eye-based communication system for users with speech impairments. *Scientific Reports* **13**, 7961 (2023).
- [8] Tarek, N. *et al.* Morse glasses: an IoT communication system based on Morse code for users with speech impairments. *Computing* **104**, 789–808 (2021).
- [9] Silva, A. B., Littlejohn, K. T., Liu, J. R., Moses, D. A. & Chang, E. F. The speech neuroprosthesis. *Nature Reviews Neuroscience* **25**, 473–492 (2024).
- [10] Card, N. S. *et al.* An Accurate and Rapidly Calibrating Speech Neuroprosthesis. *New England Journal of Medicine* **391**, 609–618 (2024).
- [11] Metzger, S. L. *et al.* A high-performance neuroprosthesis for speech decoding and avatar control. *Nature* **620**, 1–10 (2023).
- [12] Willett, F. R. *et al.* A high-performance speech neuroprosthesis. *Nature* **620**, 1031–1036 (2023).
- [13] Kim, T. *et al.* Ultrathin crystalline-silicon-based strain gauges with deep learning algorithms for silent speech interfaces. *Nature Communications* **13**, 5815 (2022).
- [14] Tang, C. *et al.* Ultrasensitive textile strain sensors redefine wearable silent speech interfaces with high machine

learning efficiency. *npj Flexible Electronics* **8**, 27 (2024).

[15] Yang, Q. *et al.* Mixed-modality speech recognition and interaction using a wearable artificial throat. *Nature Machine Intelligence* **5**, 169–180 (2023).

[16] Xu, S. *et al.* Force-induced ion generation in zwitterionic hydrogels for a sensitive silent-speech sensor. *Nature Communications* **14**, 219 (2023).

[17] Che, Z. *et al.* Speaking without vocal folds using a machine-learning-assisted wearable sensing-actuation system. *Nature Communications* **15**, 1873 (2024).

[18] Wand, M. *et al.* Tackling speaking mode varieties in EMG-based speech recognition. *IEEE Transactions on Biomedical Engineering* **61**, 2515–2526 (2014).

[19] Liu, H. *et al.* An epidermal sEMG tattoo-like patch as a new human–machine interface for patients with loss of voice. *Microsystems & Nanoengineering* **6**, 16 (2020).

[20] Wang, Y. *et al.* All-weather, natural silent speech recognition via machine-learning-assisted tattoo-like electronics. *npj Flexible Electronics* **5**, 20 (2021).

[21] Tian, H. *et al.* Bioinspired dual-channel speech recognition using graphene-based electromyographic and mechanical sensors. *Cell Reports Physical Science* **3**, 101075 (2022).

[22] Tang, C. *et al.* A deep learning-enabled smart garment for accurate and versatile sleep conditions monitoring in daily life. *arXiv.org* <https://arxiv.org/abs/2408.00753> (2024).

[23] Sherstinsky, A. Fundamentals of recurrent neural network (RNN) and long short-term memory (LSTM) network. *Physica D: Nonlinear Phenomena* **404**, 132306 (2020).

[24] Vaswani, A. *et al.* Attention is all you need. *Advances in Neural Information Processing Systems* 6000 - 6010 (2017).

[25] Chen, Z., *et al.* Long sequence time-series forecasting with deep learning: A survey. *Information Fusion* **97**, 101819 (2023).

[26] Bengio, Y., *et al.* Learning long-term dependencies with gradient descent is difficult. *IEEE transactions on neural networks* **5**, 157-166 (1994).

[27] Kiranyaz, S., *et al.* "1D convolutional neural networks and applications: A survey." *Mechanical systems and signal processing* **151**, 107398 (2021).

[28] Tang, W., *et al.* Rethinking 1d-cnn for time series classification: A stronger baseline." *arXiv preprint arXiv:2002.10061* (2020).

[29] Hinton, G. Distilling the Knowledge in a Neural Network. *arXiv preprint arXiv:1503.02531* (2015).

[30] McInnes, L., *et al.* Umap: Uniform manifold approximation and projection for dimension reduction. *arXiv preprint arXiv:1802.03426* (2018).

[31] Yu, Y., *et al.* Cloud-edge collaborative depression detection using negative emotion recognition and cross-scale facial feature analysis. *IEEE transactions on industrial informatics* **19**, 3088-3098 (2022).

[32] Yang, K., *et al.* "Behavioral and physiological signals-based deep multimodal approach for mobile emotion recognition." *IEEE Transactions on Affective Computing* **14**, 1082-1097 (2021).

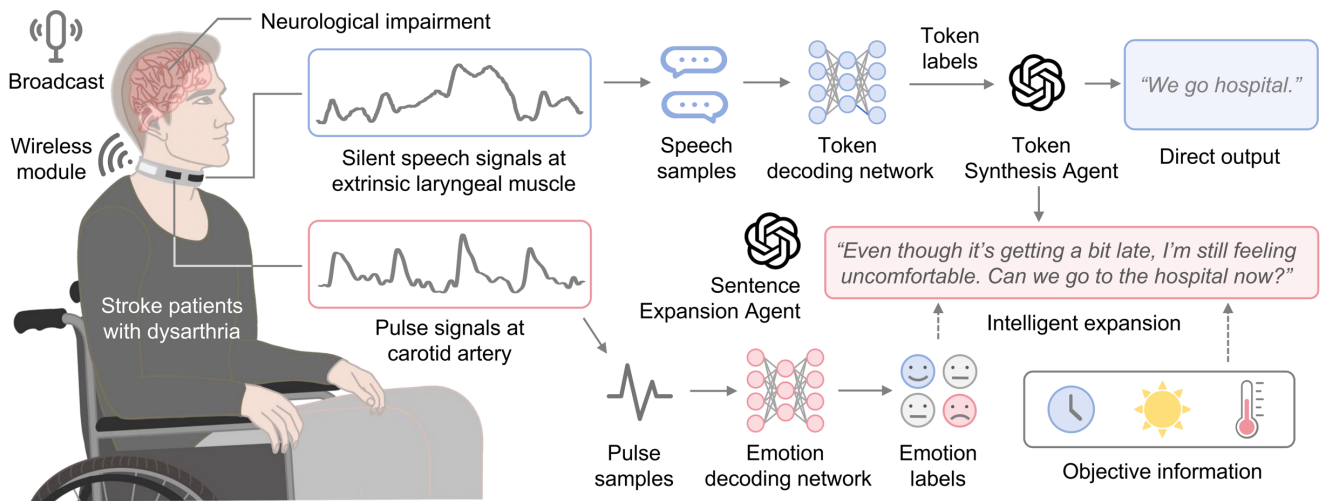
[33] Saganowski, S., *et al.* Emotion recognition for everyday life using physiological signals from wearables: A systematic literature review. *IEEE Transactions on Affective Computing* **14**, 1876-1897 (2022).

[34] Yi, W., *et al.* Ultrasensitive Textile Strain Sensing Choker for Diverse Healthcare Applications. *2024 IEEE BioSensors Conference (BioSensors)*. IEEE, 2024.

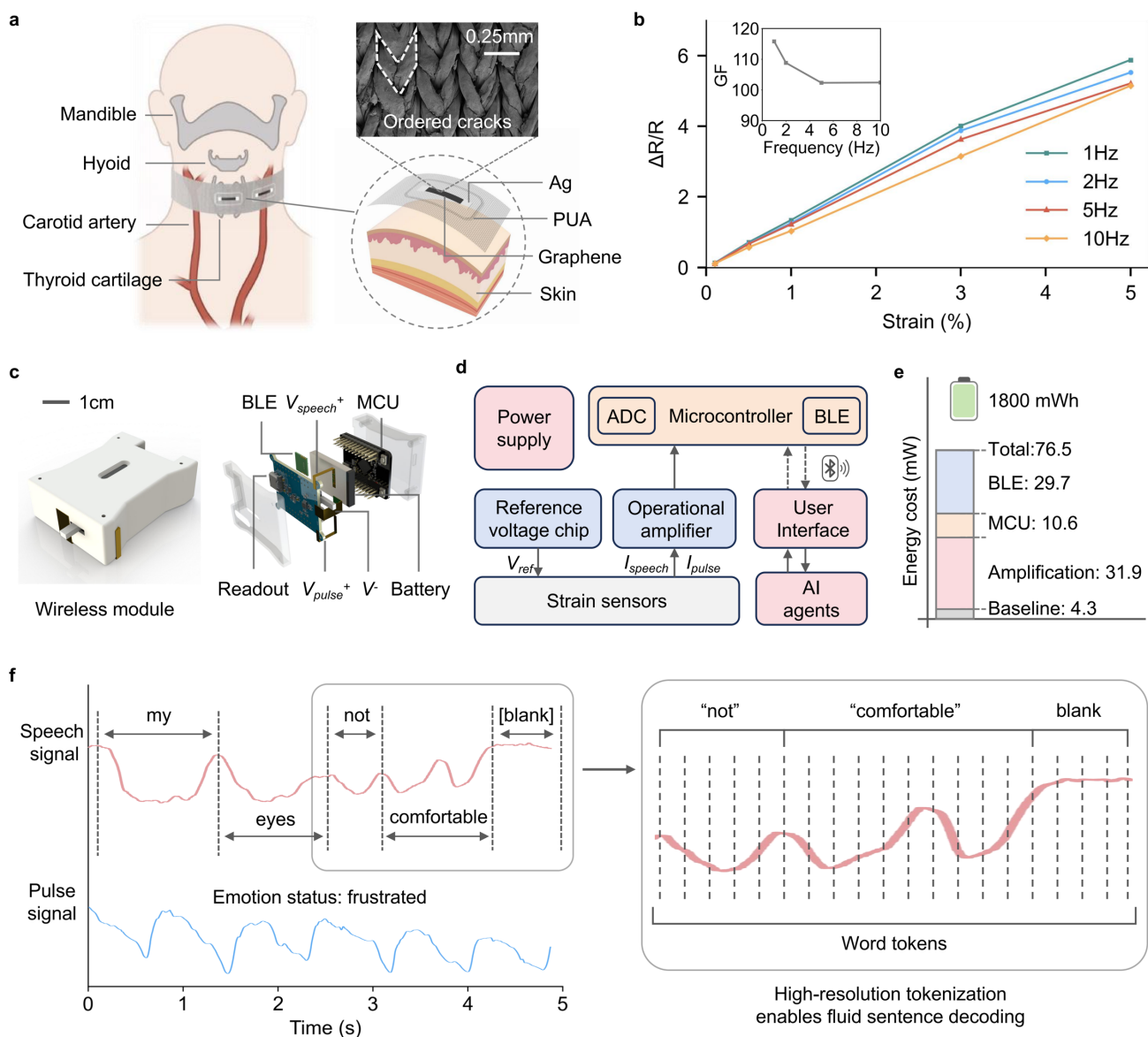
[35] Yin, J., *et al.* Motion artefact management for soft bioelectronics. *Nature Reviews Bioengineering* **2**, 541–558 (2024).

- [36] Selesnick, I., *et al.* Generalized digital Butterworth filter design. *IEEE Transactions on signal processing* **46**, 1688-1694 (1998).
- [37] Kuo, S., and Dennis M. Active noise control: a tutorial review. *Proceedings of the IEEE* **87**, 943-973 (1999).
- [38] Xie, Y., *et al.* Defending chatgpt against jailbreak attack via self-reminders. *Nature Machine Intelligence* **5**, 1486-1496 (2023).
- [39] Wei, J., *et al.* Chain-of-thought prompting elicits reasoning in large language models. *Advances in neural information processing systems* **35**, 24824-24837 (2022).
- [40] Irie, G., *et al.* Affective audio-visual words and latent topic driving model for realizing movie affective scene classification. *IEEE Transactions on Multimedia* **12**, 523-535 (2010).
- [41] Zhang, S., *et al.* Learning affective features with a hybrid deep model for audio-visual emotion recognition. *IEEE transactions on circuits and systems for video technology* **28**, 3030-3043 (2017).
- [42] Qi, Y. *et al.*, Piezoelectric Touch Sensing and Random-Forest-Based Technique for Emotion Recognition. *IEEE Transactions on Computational Social Systems* **11**, 6296-6307 (2024).
- [43] Yang, W., *et al.* Affective auditory stimulus database: An expanded version of the International Affective Digitized Sounds (IADS-E). *Behavior Research Methods* **50**, 1415-1429 (2018).
- [44] Anastassiou, P., *et al.* Seed-TTS: A Family of High-Quality Versatile Speech Generation Models. *arXiv preprint arXiv:2406.02430* (2024).
- [45] Hickok, G., and Poeppel, D. The cortical organization of speech processing. *Nature Reviews Neuroscience* **8**, 393-402 (2007).

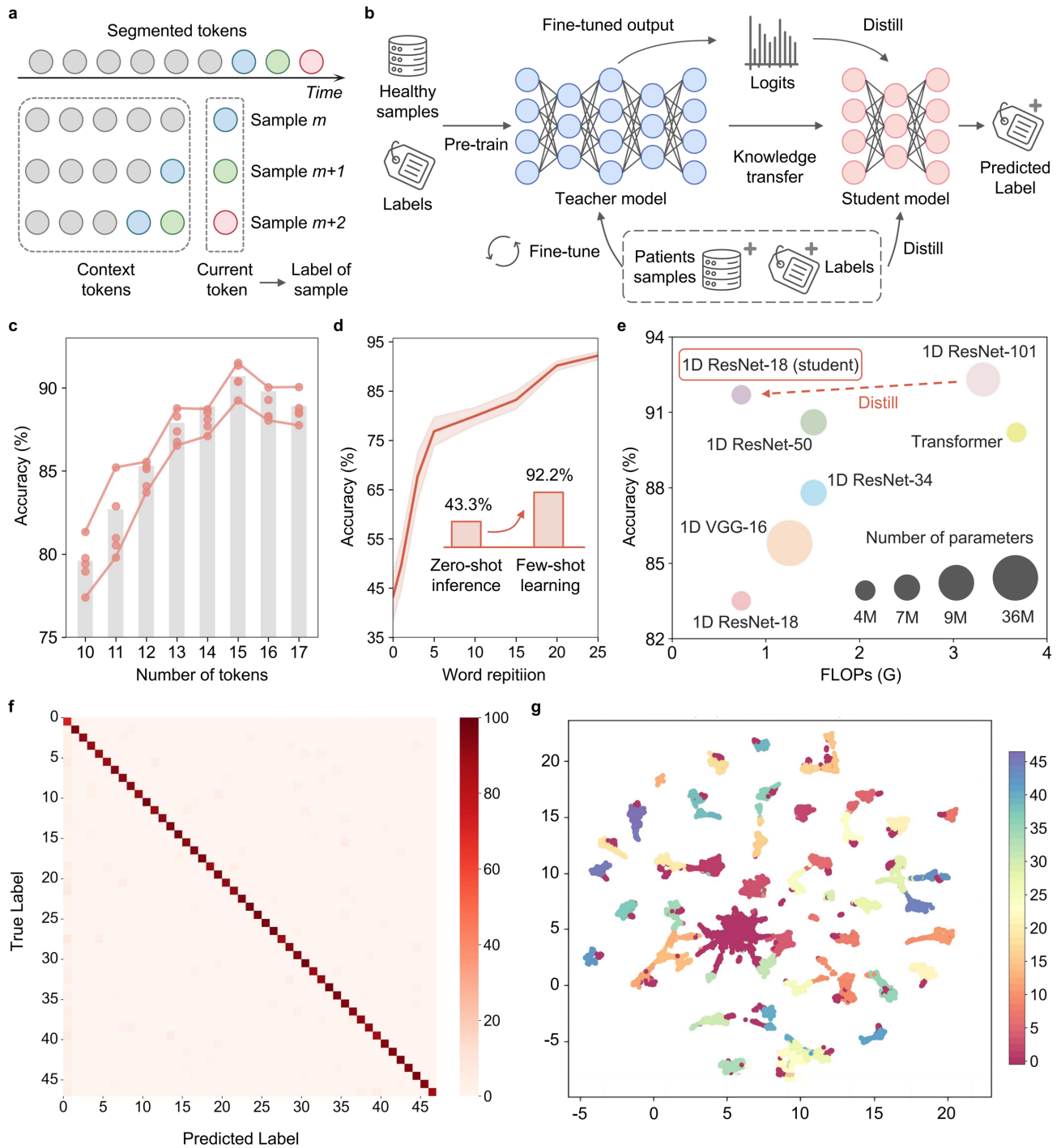
## Figures



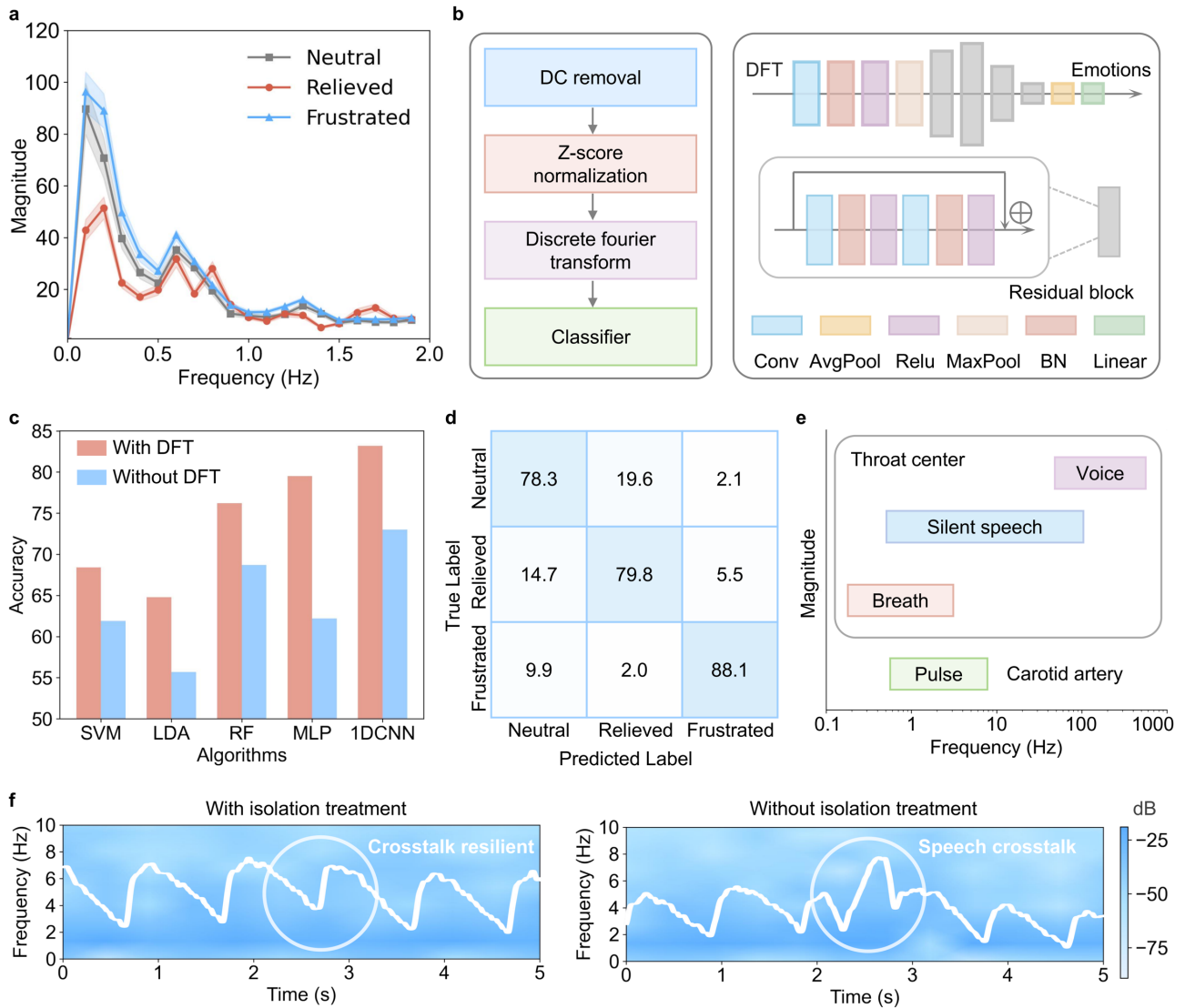
**Figure 1: Schematic of the IT developed for stroke patients with dysarthria.** The system captures extrinsic laryngeal muscle vibrations and carotid pulse signals via textile strain sensors and transmits them to the server through a wireless module. Silent speech signals are processed through a token decoding network, which generates token labels for sentence synthesis. Simultaneously, pulse signals are processed by an emotion decoding network to identify emotional states. The system intelligently integrates both emotional states and contextual objective information (e.g., time, environment) to expand the initial decoded sentences. Through a sentence expansion agent, the decoded output is transformed into personalized, fluent, and emotionally expressive sentences, enabling patients to communicate with a fluency and naturalness comparable to healthy individuals. (Note: Due to grammatical differences between Chinese and English, “We go hospital” is a word-for-word translation of the Chinese expression for “Let's go to the hospital”.)



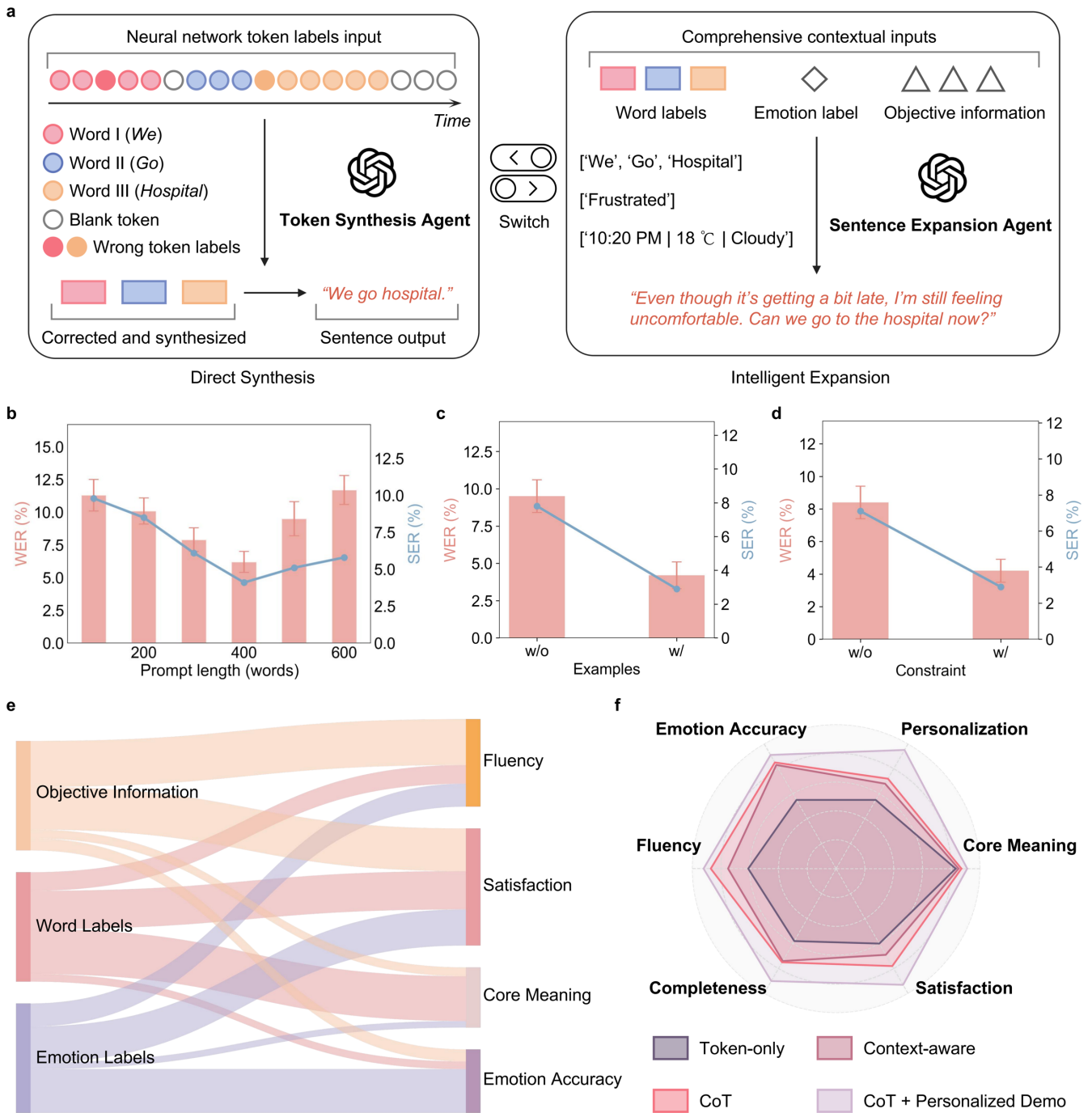
**Figure 2: Hardware and data collection of the IT.** **a**, Schematic of a textile-based strain-sensing choker. Two channels are aligned with the carotid artery and center of throat, respectively. Each channel consists of a two-terminal crack-based resistive strain sensor surrounded by a polyurethane acrylate (PUA) stress isolation layer. The top right SEM image shows the spontaneous ordered crack structure of the graphene coating. **b**, Relationship between the response to uniaxial stretching (from 0.1% to 5%) and frequency. **c**, Exploded view of the internal components of the PCB. **d**, Diagram of the system communication. **e**, Power consumption of each component during system communication. **f**, Schematic of the high-resolution tokenization strategy.



**Figure 3: Token-level decoding framework and performance evaluation.** **a**, Explicit context augmentation strategy designed to incorporate contextual information by combining tokens into token samples. **b**, Model training pipeline: the teacher model is pre-trained on healthy samples, then fine-tuned on patient samples; knowledge distillation transfers learned features to a student model for efficient prediction. **c**, Comparison of decoding accuracy across different numbers of tokens per sample, showing optimal performance when sufficient contextual information is included. **d**, Accuracy improvement with word repetition in transfer learning process, demonstrating a jump from zero-shot inference (43.3%) to few-shot learning (92.2%) as repetitions increase. **e**, Comparison of model performance across architectures with varying accuracy, FLOPs, and parameter counts; ResNet-101 and ResNet-18 were selected as the teacher and student models, respectively. **f**, Confusion matrix for the final student model. **g**, UMAP visualization of extracted features from the student model, illustrating token clustering patterns that indicate effective decoding and clear separation of different classes.



**Figure 4: Emotion decoding framework and performance evaluation.** **a**, Frequency domain characteristics of carotid pulse signals across three emotional states (Neutral, Relieved, and Frustrated), showing distinct amplitude patterns. **b**, Emotion classification workflow: preprocessing pipeline (left) involving DC removal, Z-score normalization, and discrete Fourier transform (DFT), feeding into a classifier based on a 1DCNN architecture (right) for emotion decoding. **c**, Comparison of classification accuracies across machine learning algorithms (SVM, LDA, RF, MLP, and 1DCNN) with and without DFT preprocessing, highlighting improved performance with DFT. **d**, Confusion matrix for emotion classification. **e**, Frequency and magnitude range of different vibrational signal sources (voice, silent speech, breath, carotid pulse) at neck area. **f**, Time-frequency spectrogram of pulse signals with and without strain isolation treatment when vowel “a” both introduced at 2.5s, demonstrating successful mitigation of speech crosstalk interference after applying the isolation technique.



**Figure 5: LLM agents framework and performance evaluation.** **a**, Schematic of the IT’s LLM agents: Token Synthesis Agent (left) directly synthesizes sentences from neural network token labels, while Sentence Expansion Agent (right) enhances outputs with contextual and emotional inputs. **b**, Effect of prompt length on word error rate (WER) and sentence error rate (SER) with optimal performance observed at medium lengths. **c**, Influence of example-based few-shot learning on WER and SER, showing a significant reduction when examples are provided. **d**, Impact of constrained decoding on WER and SER, demonstrating improved accuracy and sentence structure. **e**, Contribution of objective information, word, and emotion labels on key user metrics, including fluency, satisfaction, core meaning, and emotional accuracy (evaluated through ablation experiments). **f**, Radar plot comparing performance across various configurations (Token-only, Context-aware, Chain-of-Thought (CoT), and CoT with personalized demonstration) on fluency, personalization, core meaning, satisfaction, completeness, and emotion accuracy.