

Article

Not peer-reviewed version

---

# A Vision-Based Subtitle Generator: Text Reconstruction via Subtle Vibrations from Videos

---

[Yan Wang](#), Yingchong Wang, Xiuqi Zhang, [Xiaoyu Ding](#)\*

Posted Date: 11 December 2025

doi: 10.20944/preprints202512.1053.v1

Keywords: text reconstruction from vibrations; phase-based motion estimation (PME); pretrained acoustic model; transformer



Preprints.org is a free multidisciplinary platform providing preprint service that is dedicated to making early versions of research outputs permanently available and citable. Preprints posted at Preprints.org appear in Web of Science, Crossref, Google Scholar, Scilit, Europe PMC.

Copyright: This open access article is published under a [Creative Commons CC BY 4.0 license](#), which permit the free download, distribution, and reuse, provided that the author and preprint are cited in any reuse.

Disclaimer/Publisher's Note: The statements, opinions, and data contained in all publications are solely those of the individual author(s) and contributor(s) and not of MDPI and/or the editor(s). MDPI and/or the editor(s) disclaim responsibility for any injury to people or property resulting from any ideas, methods, instructions, or products referred to in the content.

Article

# A Vision-Based Subtitle Generator: Text Reconstruction via Subtle Vibrations from Videos

Yan Wang, Yingchong Wang, Xiuqi Zhang and Xiaoyu Ding \*

School of Mechanical Engineering, Beijing Institute of Technology, Haidian District, Beijing 100081, China

\* Correspondence: xiaoyu.ding@bit.edu.cn

## Abstract

Subtle vibrations induced in everyday objects by ambient sound, especially speech, carry rich acoustic cues that can potentially be transformed into meaningful text. This paper presents a Vision-based Subtitle Generator (VSG). This is the first attempt to recover text directly from high-speed videos of sound-induced object vibrations using a generative approach. To this end, VSG introduces a phase-based motion estimation (PME) technique that treats each pixel as an “independent microphone”, and extracts thousands of pseudo-acoustic signals from a single video. Meanwhile, the pretrained Hidden-unit Bidirectional Encoder Representations from Transformers (HuBERT) serves as the encoder of the proposed VSG-Transformer architecture, effectively transferring large-scale acoustic representation knowledge to the vibration-to-text task. These strategies significantly reduce reliance on large volumes of video data. Experimentally, text was generated from vibrations induced in a bag of chips by AISHELL-1 audio samples. Two VSG-Transformer variants with different parameter scales (Base and Large) achieved character error rates of 13.7 and 12.5%, respectively, demonstrating the remarkable effectiveness of the proposed generative approach. Furthermore, experiments using signal upsampling techniques showed that VSG-Transformer performance was promising even when low-frame-rate videos were used, highlighting the potential of the proposed VSG approach in scenarios featuring conventional cameras.

**Keywords:** text reconstruction from vibrations; phase-based motion estimation (PME); pretrained acoustic model; transformer

---

## 1. Introduction

When sound interacts with an object, the sound induces subtle vibrations on the object surface. The resulting motion patterns partially preserve informative features of the surrounding acoustic environment. Remarkably, if speech is a component of the ambient sound, the vibrations may encode semantic content that can be extracted and reconstructed into text. This is analogous to subtitle generation for a silent film, and we term the system that we built a Vision-based Subtitle Generator (VSG).

In recent years, subtle sound-induced vibrations of object surfaces have exhibited significant potential in terms of surveillance and security [1]. Such vibrations can reveal human activity in the environment and even enable remote eavesdropping. Most prior work in this area falls into two main categories [2]: Audio recovery and acoustic source classification.

Audio recovery methods reconstruct acoustic signals in a form that can be played and recognized by the human ear [1–9]. One classical example involves the use of laser Doppler vibrometers (LDVs) [3] that recover sound by measuring the Doppler shift in laser light reflected as the surface of an object vibrates. In recent years, researchers have explored a range of alternative physical mechanisms for audio recovery. The Lamphone [2,4] use telescopes and optical sensors to detect minute optical variations of distant objects that are induced by sound-driven vibrations. The intensities of speaker light-emitting diode (LED) indicators [5] and vibrations of the read/write heads of hard disk drives [6] can also “accidentally” leak acoustic information that enables remote audio

recovery. A fundamentally different line of research involves audio reconstruction using machine vision algorithms. This has attracted increasing attention in recent years [7,8]. Davis et al. [1] were the first to introduce a phase-based algorithm that retrieved acoustic signals from human-imperceptible motions in high-speed videos. Zhang et al. [7] developed a more computationally efficient method using singular value decomposition (SVD) techniques to enhance the robustness and accuracy of sound reconstruction from video data.

In contrast, acoustic source classification techniques are aimed at extracting information on the sound sources in a given scene. Such data include the number and/or gender of speakers. Also, the methods may link signals to particular words in a precompiled dictionary [10–14]. For example, Gyrophone [10] uses gyroscopes of the micro-electro-mechanical system embedded in modern smartphones to capture signals. This identifies speakers and enables digit-level speech recognition, albeit within a limited vocabulary (“one”, “two”, “three”, etc.). Side Eye [11] leverages the rolling shutter and movable lens mechanisms of mainstream smartphone cameras to create a point-of-view (POV), optical-acoustic side channel that accurately identifies spoken digits, speakers, and genders. As the latter technique does not focus on human auditory perception, a broad range of unintelligible or rough signals can be utilized. Table 1 presents a summary of related work.

**Table 1.** Summary of related work.

Method	Exploited Device	Sampling Rate	Technique Category
Lamphone [2,4]	Photodiode	2-4 kHz	
LDVs [3]	Laser transceiver	40 kHz	
Glowworm [5]	Photodiode	4-8 kHz	
Hard Drive of Hearing [6]	Magnetic hard drive	17 kHz	Recovery
Visual	High-speed camera	2-20 kHz	
Microphone [1]			
SVD [7,8]	High-speed camera	2.2 kHz	
SPEAKE(a)R [9]	Speakers	48 kHz	
Gyrophone [10]	Gyroscope	200 Hz	
Side Eye [11]	Smartphone cameras	60 Hz	
Accelword [12]	Accelerometer	200 Hz	Classification
PitchIn [13]	Fusion of several motion sensors	2 kHz	
WiHear [14]	Software-defined radio	300 Hz	
<b>VSG of the present paper</b>	<b>High-speed camera</b>	<b>2-16 kHz</b>	<b>Generation</b>

However, both mainstream approaches are associated with distinct challenges [2]:

(1) Audio recovery techniques focus on human auditory perception, imposing strict requirements on the complex signal processing pipelines that are often built using expert prior knowledge. Recognition performance can vary significantly by the hearing capacity and training level of the listener.

(2) The principal limitation of classification-based techniques is the restricted output space of existing models. Classification is often limited to isolated words or digits, usually in precompiled dictionaries. This renders the compilation of task-specific word lists difficult.

Innovatively, VSG represents the first attempt to recover text directly from high-speed videos of sound-induced object vibrations via a generative approach. Compared to existing constrained

approaches—either human-auditory-dependent audio recovery or fixed-vocabulary acoustic source classification—VSG transcends their limitations by formulating the task as an open-ended generative problem. To tackle Challenge 1, VSG introduces a phase-based motion estimation (PME) technique that treats each pixel as an “independent microphone”, extracting thousands of pseudo-acoustic signals (PASs) from a single video. Large volumes of human-unintelligible PASs are directly utilized during training. This also avoids any need for the complex, task-specific, audio preprocessing steps that were traditionally used to enhance the auditory experience of the listener. For Challenge 2, VSG employs an autoregressive generative architecture (termed VSG-Transformer) rather than classification-based designs. The state-of-the-art pretrained Hidden-unit Bidirectional Encoder Representations from Transformers (HuBERT) serves as the encoder in this architecture. In recent years, the original BERT [15] and variants thereof, such as HuBERT [16], XLNet [17], and RoBERTa [18], have become the dominant paradigms of natural language processing. Through transfer learning, the pretrained HuBERT model can be modified to handle downstream tasks, effectively transferring large-scale acoustic representation knowledge to the vibration-to-text conversion. Moreover, these strategies significantly reduce reliance on large volumes of video data.

Experimentally, the vibrational responses of everyday objects are converted to standard audio samples of the AISHELL-1 corpus. This is a widely used, open-source Mandarin speech dataset that serves as a speech recognition benchmark [19]. The evaluation metric was the character error rate (CER). This is the standard parameter employed when assessing automatic speech recognition (ASR) models. The results validate the effectiveness of the proposed VSG approach. Text is reconstructed from subtle object vibrations captured on video. Specifically, VSG-Transformer-Base and VSG-Transformer-Large (similar models but of different scales) exhibited CERs of 13.7 and 12.5%, respectively. The applicability of VSG-Transformer to low-frame-rate videos was further investigated. Signal upsampling techniques were used to alleviate any dependence on high-speed imaging devices. When the interpolation strategies were appropriate, VSG-Transformer maintained an acceptable recognition performance even under the usual low-frame-rate constraints. This means that the VSG approach can be used even when conventional cameras collect scene data.

The remainder of the paper is organized as follows: Section 2 describes the proposed method. Section 3 details the experimental validation and discusses the results. Section 4 explores whether lightweight VSG implementation is possible when low-frame-rate videos are processed using upsampling techniques. The conclusions are in Section 5.

## 2. Methods

This section presents the general framework of the proposed VSG and the key related concepts, including PAS synthesis and the VSG-Transformer architecture. The theoretical background and implementation details follow.

### 2.1. General Framework of the Method

VSG evaluates only ambient sound. The semantic content encoded in pixel signals is used to reconstruct output text. Figure 1 illustrates the overall flow, which consists of the following steps (A full demonstration animation of VSG is available at the provided [[Link 1](#)] for readers):

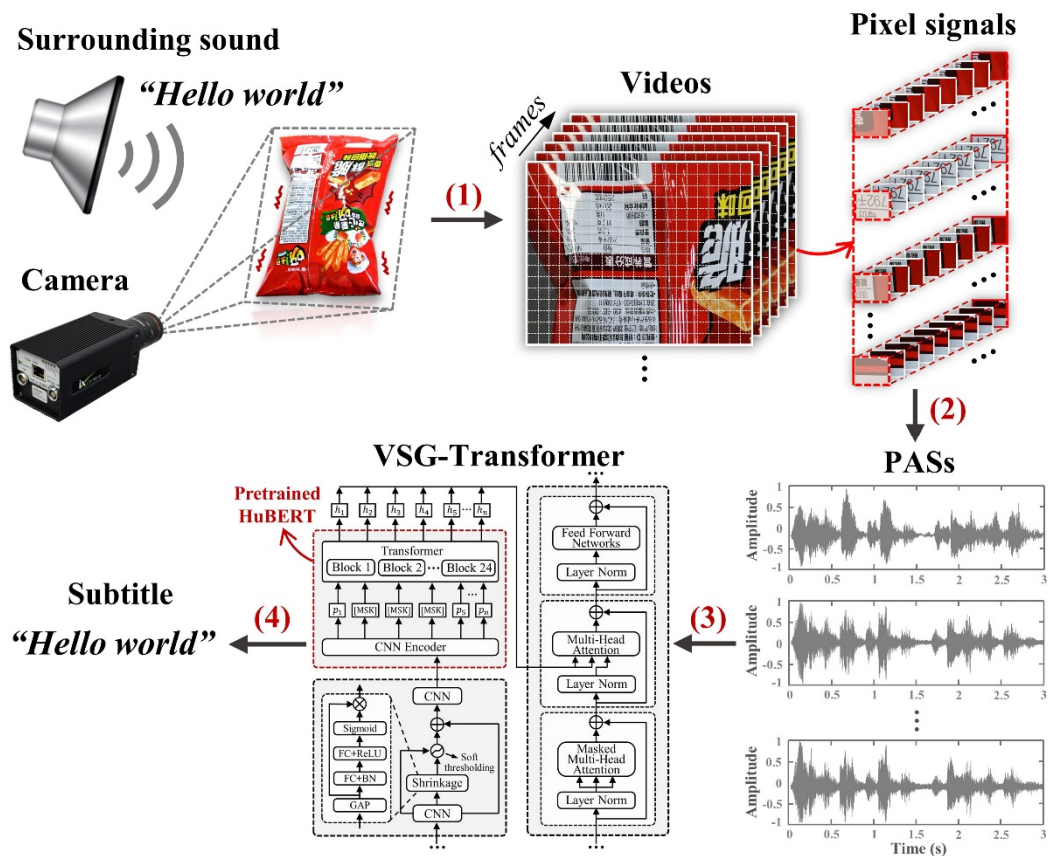


Figure 1. Flowchart of the proposed VSG method.

1. Video capture: The response of any object to sound is purely physical. As the acoustic excitations vary, the resulting object surface vibrations are captured by a camera, effectively transforming physical displacements into pixel-level signals within the video frames;
2. PAS acquisition: PASs are obtained via phase-based processing of the pixel signals;
3. VSG-Transformer training and testing: Large-scale PASs that encode rich acoustic features are used to construct the PAS dataset employed to train and evaluate VSG-Transformer. A multi-stage transfer learning strategy effectively links the pretrained acoustic representations of HuBERT to the PAS-driven VSG task;
4. Text reconstruction: The trained VSG-Transformer reconstructs text based the PASs extracted from new videos.

## 2.2. Extraction of PAS

As outlined above, VSG requires audio-like signals that effectively encode acoustic features of the input video  $V(x,y,t)$ . Crucially, such signals must be widely accessible. This section details the principles and procedures involved. The process can be broadly divided into two stages. First, local motion signals are computed at each pixel location. Next, these signals are transformed into pixel-level representations that approximate audio waveforms. These are the PASs.

### (a) Computation of Local Motion Signals

Fleet and Jepson were the first to use spatio-temporal bandpass filters for PME in image sequences [20]. The phase gradient of the complex-valued output is a reliable approximation of the motion field. Building on this, subsequent studies used more advanced filtering strategies [21–23]. For example, Gautama and Van Hulle [21] employed a set of quadrature Gabor filter pairs when computing the temporal phase gradient of a spatially bandpassed video to estimate the motion field. In recent years, many studies have bypassed the explicit computation of optical-flow vector fields by,

rather, directly leveraging phase variations when estimating the displacements of image textures in video sequences [1,24].

Following the approach of [1], local motion is here computed by analyzing the phase variations within a complex, steerable pyramid representation of the input video  $V(x,y,t)$ . A complex steerable pyramid is a multi-scale, multi-orientation filter bank (see [24] for details) that decomposes each video frame into complex-valued sub-bands, indexed by the scale  $r$  and orientation  $\theta$ . At each pixel location  $(x,y)$ , the sub-band output can be represented in terms of amplitude  $A$  and phase  $\varphi$ , as follows:

$$A(r, \theta; x, y, t) e^{i\varphi(r, \theta; x, y, t)} \quad (1)$$

Traditionally, a complex steerable pyramid is subjected to downsampling-based decomposition as in [1,24]. Here, however, the resolution across all filters is uniform (Figure 2). This yields sub-band outputs that are perfectly aligned in terms of spatial resolution. The amplitude and phase at each pixel location  $(x,y)$  across the different sub-bands exhibit a direct correspondence. This facilitates later pixel-level signal synthesis.

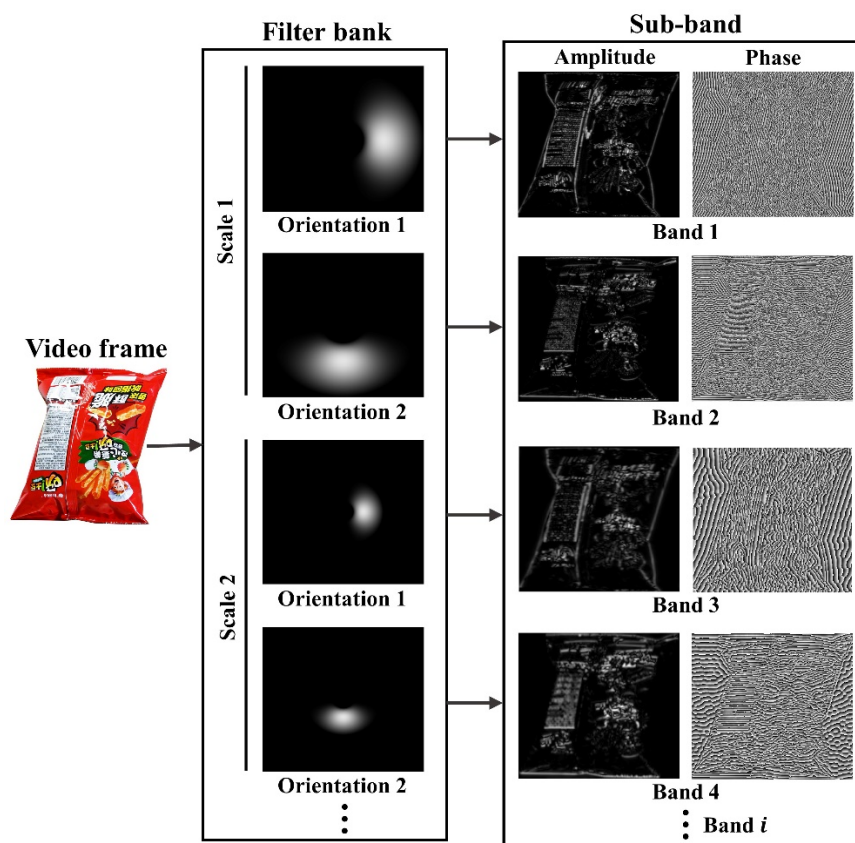


Figure 2. Decomposition of a video frame.

The phase variations  $\varphi_v(r, \theta; x, y, t)$  are then computed by subtracting the phase of each pixel in the reference frame—typically the first frame of the video—from the corresponding phase in the current frame. The formal expression is:

$$\varphi_v(r, \theta; x, y, t) = \varphi(r, \theta; x, y, t) - \varphi_{ref}(r, \theta; x, y, t_0) \quad (2)$$

By [21], the computed phase variations afford very good approximations of image texture displacements, especially those of subtle motions, along the corresponding orientation and scale.

(b) PAS Synthesis

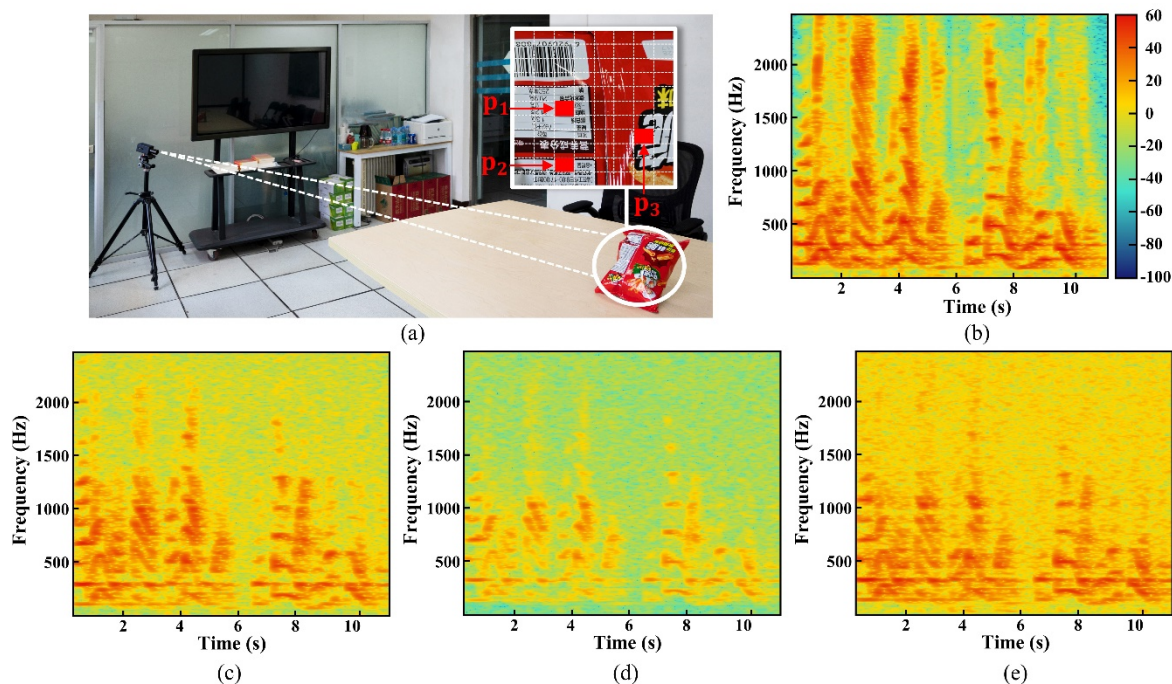
A large number of PASs can be independently constructed by utilizing phase variations  $\varphi_v(r, \theta; x, y, t)$  at each pixel. Specifically, within each sub-band, each local motion signal is weighted by the squared amplitude at time  $t_0$ . This can be formulated as:

$$\Phi_i(r, \theta; x, y, t) = A(r, \theta; x, y, t_0)^2 \varphi_v(r, \theta; x, y, t) \quad (3)$$

where  $i$  denotes the sub-band index. Subsequently, the outputs across all orientations  $\theta$  and scales  $r$  are aggregated via summation:

$$PAS(x, y, t) = \sum_i \Phi_i(r, \theta; x, y, t) \quad (4)$$

Finally, PAS is normalized by scaling and centering to within the range  $[-1, 1]$ . Figure 3 compares the spectrogram of the original audio [Figure 3(b)] to those of the PASs [Figure 3(c), (d), and (e)] derived from different surface regions [ $p_1, p_2, p_3$  of Figure 3(a)] of the object. To provide an intuitive understanding of signal quality, the PASs were rendered as audible waveforms using the method proposed in [1]. Audio samples corresponding to Figure 3(b–e) are also available at the provided [Link 2], enabling subjective comparison of the reconstructed PASs with the original audio. Given the inherent ambiguities of local phases in regions where image texture is weak, motion signals extracted from such pixel locations tend to be very noisy and/or unreliable.



**Figure 3.** A comparison between the spectrogram of the original audio (b) and those of PASs [(c), (d), (e)] obtained from surface locations  $p_1, p_2, p_3$ . The PASs were extracted from a 4.4-kHz video of a bag of potato chips on the table of a typical meeting room [Link 2]. The sentence “We are continuously improving our technology”, corresponding to the Mandarin pronunciation “Wǒmen yězài bùduàn tǐshēng wǒmen de jìshù,” was spoken by a person near the bag at an approximate volume of 80 dB.

In most areas, the resulting PASs are too coarse to be intelligible to the human ear. Current end-to-end ASR models also exhibit limited recognition performance on these signals, as detailed in Section 3.3. However, the advantage of PASs lies in their abundance—they can be densely extracted at the pixel level and collectively capture rich semantic information embedded in surface vibrations. Experimental results confirm their substantial training potential when combined with dedicated network architectures and learning strategies, as further demonstrated in Section 4.

Note that variations in texture, brightness, and material properties across different pixel locations that correspond to distinct regions of the object surface introduce diverse noise

characteristics and frequency attenuations. Such variations contribute to the robustness of model training using large-scale PAS inputs.

### 2.3. Proposed Model: The VSG-Transformer

Figure 4 shows the architecture of the VSG-Transformer. An overview of each component follows, and the design rationale is explained. In broad terms, the model features three key modules: (1) a convolutional shrinkage module that suppresses noise and enhances features; (2) a pretrained HuBERT-based encoder that captures high-level acoustic representations; and, (3) an attention-based decoder that generates autoregressive text.

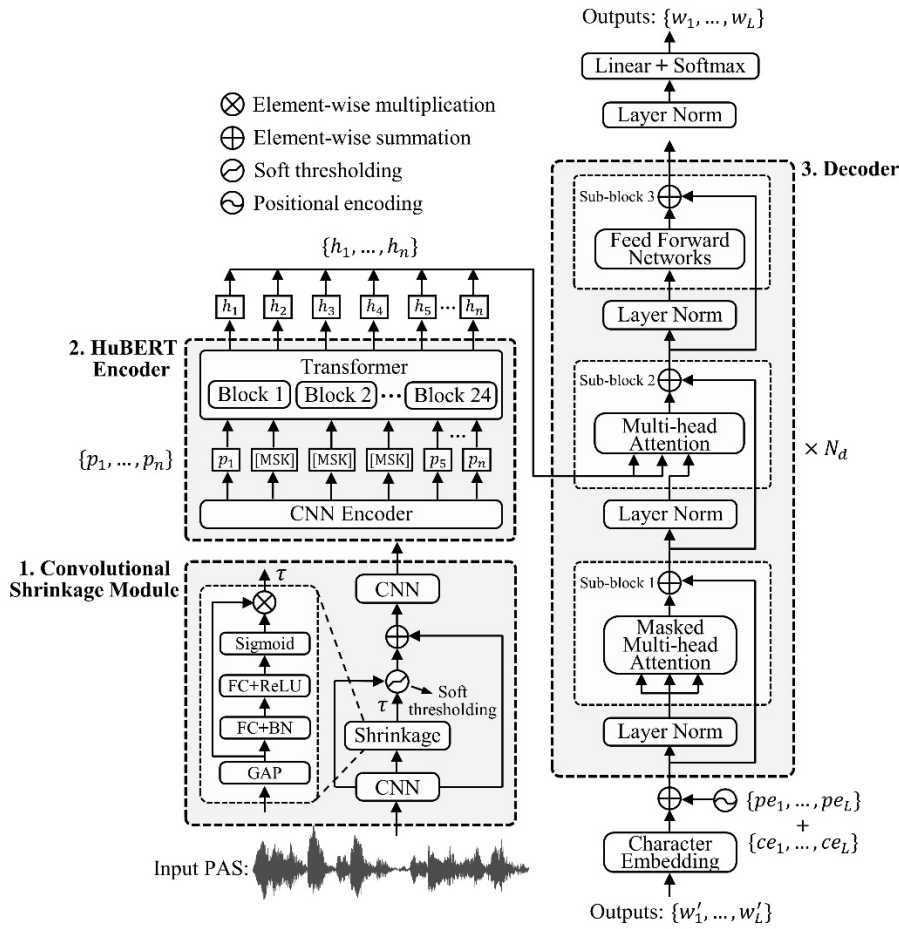


Figure 4. The VSG-Transformer architecture.

#### (a) The convolutional shrinkage module

As defined by Equation (4) of Section 2.2, PASs are combinations of  $\Phi_i$  derived from the amplitude and phase outputs across all orientations  $\theta$  and scales  $r$  of the corresponding filters. For a given pixel location, filters that are mismatched in terms of scale and/or orientation typically deliver weak sub-band amplitudes. If such amplitudes are used to compute squared weightings, the PASs contain noise components that are substantially weaker than the actual signal. Therefore, soft thresholding was here used to eliminate noise. Such thresholding is as a core feature of signal denoising algorithms [25]. Noise features with magnitudes near zero are suppressed, effectively reducing them to zero. The nonlinear nature of the transformation is formally defined by the following equation:

$$y = \begin{cases} x - \tau & x > \tau \\ 0 & -\tau \leq x \leq \tau \\ x + \tau & x < -\tau \end{cases} \quad (5)$$

where  $x$  denotes the input feature,  $y$  is the output, and  $\tau$  is a positive threshold parameter. However, determination of an appropriate threshold  $\tau$  remains challenging in practice.

As shown in [26], soft thresholding can be integrated into the deep architecture as a nonlinear transformation layer. The threshold parameter  $\tau$  is then automatically learned by the network, not manually specified. Here, soft thresholding was employed as the activation function that was integrated with convolutional layers to form a convolutional shrinkage module that suppressed noise-related information and enhanced feature discriminability.

Specifically, given a PAS dataset for a VSG task with  $N$  training samples,  $D_{VSG} = \langle PAS^{(j)}, W^{(j)} \rangle_{j=1}^N$ , where PAS is a one-dimensional waveform input and  $W$  a sequence of  $L$  tokens  $\{w'_1, \dots, w'_L\}$  (Figure 4). The convolutional shrinkage module features two 1D convolutional layers that flank a shrinkage unit. The kernel configurations are 32 and 1 respectively, and both maintain a stride of 1. Initial processing through the first convolutional layer expands the feature channels of the waveform input to 32. The subsequent shrinkage layer first dynamically computes channel-specific  $\tau$  values and generates a 32-channel threshold vector. Next, that layer subjects each feature channel to a nonlinear soft thresholding operation, effectively suppressing noise but preserving critical signal components. The residual connection integrates the features that are soft-thresholded. The final 1D convolutional layer compresses the feature channels back to a single channel. This ensures that the output is compatible with the subsequent processing stages.

#### (b) The pretrained HuBERT encoder

In the second module, the pretrained HuBERT encoder extracts high-level acoustic representations from the input waveform. HuBERT has consistently attracted considerable research attention. HuBERT captures complex temporal dependencies and semantic features effectively because HuBERT was trained using large datasets, including 10000 hours of speech from the WenetSpeech [27]. The pretrained HuBERT model uses its extensive prior knowledge to engage in transfer learning. HuBERT therefore rapidly adapts to PAS inputs. In this study, two pretrained HuBERT models with different parameter scales—HuBERT-Base (95 million parameters) and HuBERT-Large (317 million parameters)—were integrated with VSG-Transformer, yielding VSG-Transformer-Base and VSG-Transformer-Large respectively.

As illustrated in Figure 4, the pretrained HuBERT model features a convolutional neural network (CNN) encoder followed by a Transformer. The CNN encoder contains seven convolutional layers, each with 512 output channels. The first layer receives a single-channel input that directly matches the one-dimensional waveform output by the convolutional shrinkage module. The subsequent layers process the 512 channel inputs. The CNN encoder delivers a downsampled two-dimensional feature sequence  $\{p_1, \dots, p_n\}$  wherein each feature vector  $p_i$  is of dimensionality 512 and is computed by aggregating information across all input channels. The feature sequence is subjected to local masking using the strategy described in [28,29] before transmission to the Transformer module, the Transformer features a configurable number of blocks with predefined embedding dimensions, inner FFN dimensions, and attention heads. The detailed architectural specifications of the HuBERT-Base and HuBERT-Large models are listed in Table 2.

**Table 2.** The parameters of HuBERT-Base and HuBERT-Large.

		Base	Large
CNN encoder	Strides	5, 2, 2, 2, 2, 2, 2	
	Kernel Width	10, 3, 3, 3, 3, 2, 2	
	Channels	512	
Transformer	Blocks	12	24
	Embedding Dimension	768	1,024
	Inner FFN Dimension	3,072	4,096

	Attention Heads	12	16
	Number of Parameters	95 M	317 M

Both models output a sequence of hidden units  $\{h_1, \dots, h_n\}$  that serve as high-level acoustic representations. These are next fed to the downstream decoder module to improve text generation.

### (c) The decoder

The decoder module of VSG-Transformer transforms the high-level acoustic representations produced by the HuBERT encoder into coherent text outputs. The decoder employs a self-attention mechanism to effectively capture long-range dependencies within the sequence. The decoder operates in an autoregressive manner through a stack of attention blocks, integrating both acoustic features from the HuBERT encoder and text-level information from previous decoding steps to enhance transcriptional accuracy.

The structure of the decoder is illustrated in Figure 4. A character-embedding layer is first used to convert the character sequence into an output encoding  $\{ce_1, \dots, ce_L\}$ , which is then combined with a positional encoding  $\{pe_1, \dots, pe_L\}$ . The dimensionality of both embeddings is that of HuBERT. The combined embeddings are then passed through a stack of  $N_d = 6$  decoder blocks that generate the final decoder outputs. Each decoder block features three sub-blocks, of which the first is a masked multi-head self-attention block wherein the queries, keys, and values are identical. Masking is used to ensure that the prediction at position  $j$  depends only on the outputs at positions earlier than  $j$ . The second sub-block is a multi-head attention block in which the keys and values are the encoder outputs  $\{h_1, \dots, h_n\}$ . Here, the queries are derived from the output of the preceding sub-block. The third sub-block is a position-wise feed-forward network. Each sub-block in the decoder features a residual connection and layer normalization. The latter uses a pre-norm rather than the conventional post-norm structure. Given input  $x$  to a sub-block, the corresponding output is:

$$x + \text{SubBlock}(\text{LayerNorm}(x)) \quad (6)$$

Finally, the decoder outputs are computed via a linear projection followed by a softmax function.

In general, VSG-Transformer effectively integrates noise suppression, acoustic representation, and autoregressive text generation, leveraging transfer learning to fully exploit the capabilities of HuBERT and enable reliable PAS-to-text conversion. The combination of modules described above is well-suited to the VSG task, which demands accurate modeling of complex acoustic and semantic patterns embedded in input PAS signals.

Notably, VSG-Transformer does not rely heavily on semantic features within PASs per se. Rather, transfer learning enables the model to leverage knowledge acquired in the source domain (the speech corpora) using the HuBERT module and then apply that information to tasks in the target domain (the PAS dataset). This significantly reduces model dependency on large volumes of PAS data. The model can be trained using only a few high-speed videos. Data acquisition from such videos is computationally demanding.

## 3. Experimental Validation

The VSG technique was evaluated in a series of experiments. The setup included a common object (a bag of potato chips), a loudspeaker (KRK Rokit 5), and a high-speed camera (i-SPEED 230). All video recordings were captured in a typical meeting room (Figure 5). Videos were captured at high frame rates (about 16 kHz; the spatial resolution was  $128 \times 128$  pixels). Standard audio files from the AISHELL-1 corpus, an open-source Mandarin speech dataset with a speech recognition baseline, were played through a loudspeaker at approximately 80 dB corresponding to a sound pressure level (SPL) comparable to that of an unamplified stage actor. A total of 500 videos were collected. Experimental evidence suggests that this quantity is adequate, as the model's semantic understanding primarily originates from the pretrained module. The PASs extracted from the videos

expose the model to variations in texture, brightness-related noise, and frequency attenuation—factors that enhance robustness during training. Subsequently, initial and final non-informative frames were manually truncated to retain only valid segments.

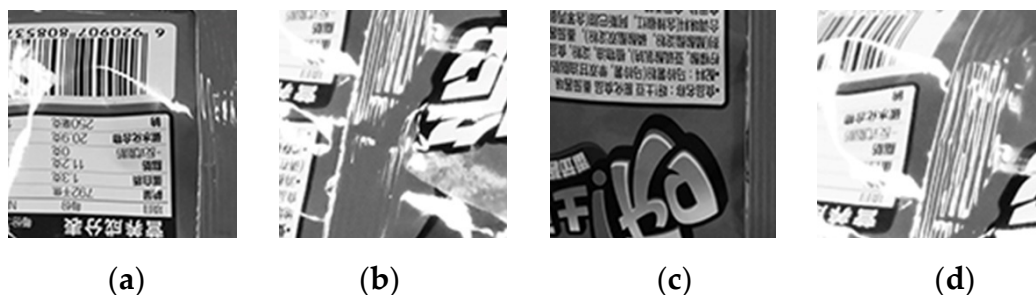


**Figure 5.** The experimental setup.

A chip bag was selected as the primary object for evaluation due to its strong and consistent vibration response under acoustic excitation. Compared to other common materials such as foam cups, plants, tissue boxes, and other everyday items, it exhibits more favorable motion characteristics, as demonstrated in [1]. Since the current experiment represents an initial feasibility study of VSG, the evaluation focused on objects with a higher likelihood of producing reliable results. Given the substantial time and computational demands of high-speed video acquisition and processing, large-scale testing across a broader range of materials is being conducted in follow-up work and will be published in future updates.

### 3.1. Dataset Generation

Figure 6 shows frames from videos captured during the experiments. Each frame is fully occupied by the object, effectively eliminating background interference. During recording, the object pose and the camera viewing angle were randomly varied but with maintenance of an approximately constant distance between the object and the camera. This significantly enhanced the robustness and generalizability of the experimental results.



**Figure 6.** Frames from the videos.

PASs were extracted from the videos as described in Section 2.2. Motion signals from pixels on the object surface that exhibited clear textural and phase variations were employed. Specifically, all PASs within each video—128 pixels (height)  $\times$  128 pixels (width) = 16384 pixels in total—were extracted for subsequent processing. A percentile-based pooling strategy was employed. This is analogous to the pooling layers employed during deep learning. Specifically, each input frame of size

128 (height)  $\times$  128 (width) pixels was partitioned into 256 regions, each measuring 8 (height)  $\times$  8 (width) pixels. Within each region, the PAS corresponding to the pixel, the amplitude of which was closest to  $A_{min} + t \times (A_{max} - A_{min})$ , was selected, where  $A_{max}$  and  $A_{min}$  were the maximum and minimum amplitude within the region, respectively, and  $t$  a user-defined threshold (here, 0.8). This identified a representative PAS within each local region, suppressed any effects of noise and outliers, and thereby improved the general quality of the selected PASs.

After application of the percentile-based pooling strategy, each 128 (height)-  $\times$  128 (width)-pixel video frame yielded 256 PASs. A total of 500 videos were randomly split into training, validation, and test sets in a 70:15:15 ratio (350, 75, and 75 videos respectively). In total, the PAS dataset contained 128000 samples (500  $\times$  256). The details are listed in Table 3.

**Table 3.** Statistics: The AISHELL-1 and PAS datasets.

		AISHELL-1			PAS Dataset		
		Train.	Dev.	Test	Train.	Dev.	Test
Utterances		120098	14326	7176	89600	19200	19200
Hours		150	18	10	107	28	29
Durations (s)	Min.	1.2	1.6	1.9	3.5	3.8	3.5
	Max.	14.5	12.5	14.7	12.4	10.2	10.4
	Avg.	4.5	4.5	5	4.3	5.3	5.5
Tokens	Min.	1.0	3.0	3.0	4.0	5.0	4.0
	Max.	44.0	35.0	37.0	35.0	22.0	26.0
	Avg.	14.4	14.3	14.6	13.6	12.3	11.2

Additionally, the publicly available AISHELL-1 corpus, which contains 178 hours of Mandarin speech, was employed during certain stages of training. All recordings were sampled at 16 kHz. To facilitate comparisons, each training sample in the PAS dataset is here termed an “utterance”, following the AISHELL-1 convention. However, as detailed in Section 2.2, PASs are not equivalent to semantically intelligible speech in the traditional sense.

### 3.2. Training and Testing

Section 2.3 described the architecture of the proposed framework. There are three key modules: a convolutional shrinkage module for noise suppression, a pretrained HuBERT encoder for acoustic representation, and an attention-based decoder for text generation. A transfer learning strategy was used to retain as much as possible of the knowledge of the pretrained HuBERT when effectively training our model across all components. This involved stage-wise adjustments of the training parameters, the use of module freezing policies, and employment of different dataset inputs. This section outlines these in detail.

During training stage 1, both the shrinkage layer of the convolutional shrinkage module and the entire HuBERT encoder were frozen, as illustrated in Figure 7(a). During this stage, only AISHELL-1 was used for training over 130 epochs. Two parallel experiments were conducted using HuBERT-Base and HuBERT-Large. Stage 1 sought to establish an automatic baseline ASR capability based on the AISHELL-1 corpus. In stage 2, the shrinkage layer was thawed but the HuBERT encoder remained frozen [Figure 7(b)]. Training in this stage employed only the PAS dataset (40 epochs). The two parallel training tracks were based on the respective models obtained during stage 1. Stage 2 aimed to optimize the shrinkage layer in terms of task-specific noise suppression while allowing transfer of knowledge from AISHELL-1 to the recognition tasks posed by the PAS dataset.

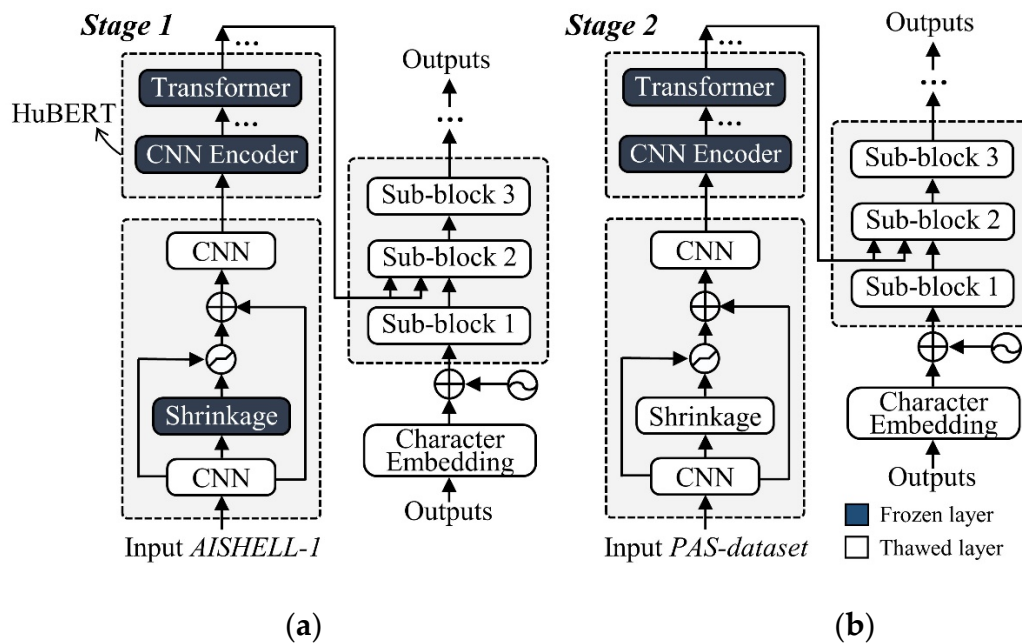


Figure 7. The stage-wise training protocol.

All VSG-Transformer modules were implemented in PyTorch [30]. Each training batch (size=24) contained approximately 100 s of speech or PAS and the corresponding transcriptions. The Adam optimizer [31] was employed; the parameters were  $\beta_1=0.9$ ,  $\beta_2=0.98$ , and  $\epsilon=10^{-9}$ . The learning rate was varied during training using the warm-up schedule of:

$$\text{lrate} = D_m^{-0.5} \cdot \min(\text{step}^{-0.5}, \text{step} \cdot \text{warmup}^{-1.5}) \quad (7)$$

where warmup was set to 12000. To prevent overfitting, the neighborhood smoothing strategy of [32] was employed, with the probability of the correct label set to 0.8. The token vocabulary included 4231 characters from the training set with two special symbols: “<unk>” for unseen tokens and “<eos>” to pad the ends of token sequences. Both the residual dropout and attention dropout rates were 0.1. Residual dropout was applied to each sub-block prior to addition of the residual connections. Attention dropout was applied to the softmax activation within each attention mechanism. Finally, the model parameters from the last 10 epochs were averaged to obtain the final output. During inference, beam search decoding employed a beam width of 10 and a length penalty of 1.0 [33]. When PASs were extracted from the same video, a majority voting strategy was used to aggregate the recognition results. This featured token-wise voting across all sequences followed by selection of the most frequent token at each position when creating the final transcription. All experiments were conducted on an NVIDIA RTX 4090 GPU, and the reported CER results are the averages over five independent runs. The overall pipeline of VSG is visualized in Figure 8.

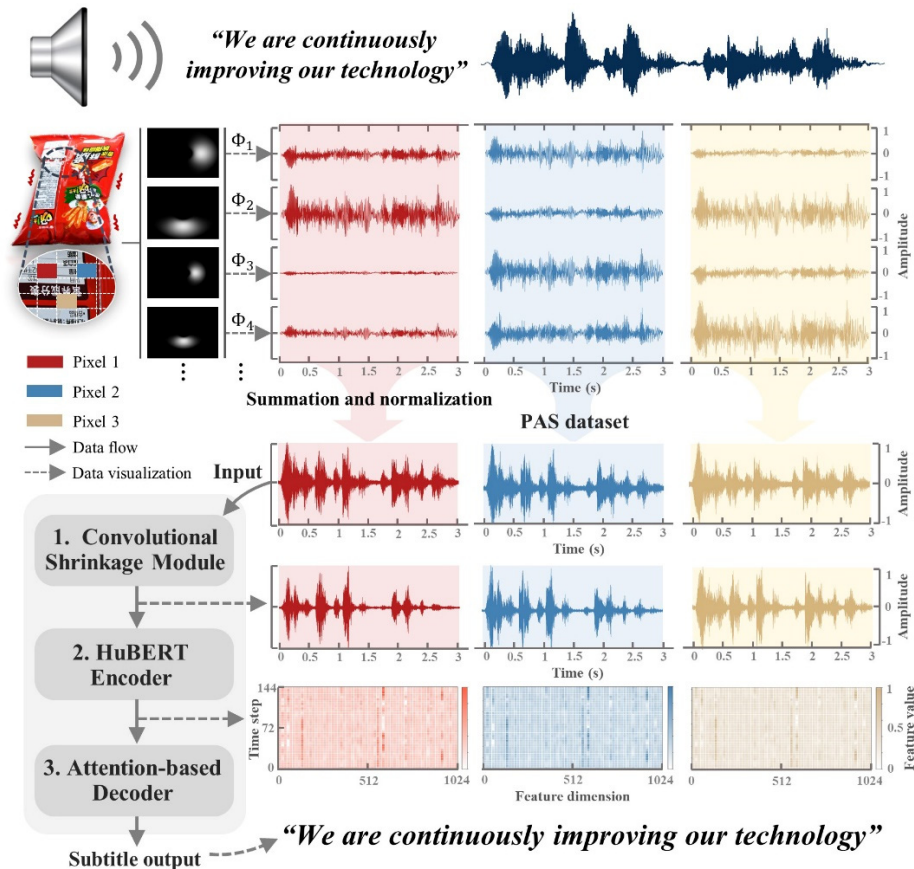


Figure 8. Visualization of the proposed VSG method.

### 3.3. Results and Ablation Studies

The results are presented in Table 4. Using the transfer learning strategy, the VSG-Transformer successfully achieved PAS-to-text conversion using high-speed video recordings of common vibrating objects. Human auditory perception was not in play. Specifically, VSG-Transformer-Base achieved a CER of 13.7% on the VSG task; the VSG-Transformer-Large figure was 12.5%. In contrast, for the AISHELL-1 ASR task of stage 1, the CERs were 6.4% and 6.1%, respectively, indicating that VSG-Transformer can engage in end-to-end ASR. VSG-Transformer-Large performed better than VSG-Transformer-Base on both the ASR and VSG tasks, showing that an increase in the model scale enhanced the capacity to extract rich acoustic representations for downstream tasks.

Table 4. Model CERs on ASR and VSG tasks at different training stages.

Training Stage	Model Scale	Dataset		Training Epochs	Frozen Layer(s)	Development (%)	Test (%)
		AISHELL-1	PAS				
Stage 1	Base	√	-	130	Shrinkage + HuBERT	6.2	6.4
	Large	√	-	130	Shrinkage + HuBERT	5.9	6.1
Stage 2	Base	-	√	40	HuBERT	13.3	13.7
	Large	-	√	40	HuBERT	12.1	12.5

After multi-stage training of the multi-module VSG-Transformer network, ablation studies were used to assess quantitatively the utilities of individual modules within the network architecture and to isolate and evaluate the impacts of the stage-wise training strategy by selectively removing certain training components.

The ablation study featured three tests. In test 1, stage 2 training was omitted, and the VSG-Transformer was therefore trained on only the AISHELL-1 dataset during stage 1. The goal of test 1 was to evaluate whether PAS signals could be directly recognized by ASR models, and to quantify their degradation relative to real audio in terms of quality and intelligibility. The model was then used to evaluate the PAS dataset test set constructed using AISHELL-1 as described in Section 3.1. The utterances of both datasets would be expected to exhibit high-level acoustic and semantic similarity. Therefore, test 1 quantitatively assessed the differences between the AISHELL-1 and PAS dataset. Test 2 preserved both stages 1 and 2 of training, but froze both the stage 2 shrinkage layer and the HuBERT encoder to assess quantitatively the impact of the shrinkage layer used for noise suppression on the final recognition accuracy. Test 3 ablated the network architecture. The performance of VSG-Transformer was examined when the number of decoder blocks ( $N_d$ ) was 4, 6, 8, or 10. All other experimental settings were those of Section 2.3. All ablation tests used the VSG-Transformer-Large model, which exhibited superior baseline performance. The results are presented in Table 5.

As shown in Table 5, the “Without stage 2” setting of test 1 was associated with poor results. The VSG-Transformer trained solely on AISHELL-1 did not perform well on the PAS-dataset test set; the CER was 56.7%. This confirms that, despite the presence of acoustic and semantic similarities between the utterances of the two datasets, domain-specific variations remain and must be addressed during stage 2 training to ensure satisfactory performance. Test 2 (“Shrinkage layer frozen during stage 2”) yielded a CER of 19.1%, indicating significant performance degradation compared to that of the full model (CER 12.5%). This emphasizes the need for soft-thresholding functions within the convolutional shrinkage module; such functions suppress noise in the PAS dataset. In test 3, the VSG-Transformer-Large model used different numbers of decoder blocks. Models with more blocks generally performed better, suggesting that a deeper decoder extracts more discriminative, token-level semantic representations. However, the performance improvements at higher decoder block numbers were rather limited, especially when the number of blocks exceeded 6.

**Table 5.** Ablation test configuration and the CER results.

	Configuration	AISHELL-1		PAS Dataset	
		Development (%)	Test (%)	Development (%)	Test (%)
Test1	Without stage 2	5.9	6.1	-	56.7
Test2	Shrinkage layer frozen during stage 2	5.9	6.1	18.5	19.1
	Number of decoder blocks	Development (%)	Test (%)	Development (%)	Test (%)
	10	5.8	6.1	12.0	12.1
Test3	8	5.8	6.1	12.0	12.2
	6 (baseline model)	5.9	6.1	12.1	12.5
	4	6.4	6.7	13.4	13.6

#### 4. VSG Operation with Low-Frame-Rate Videos

The VSG framework and the experimental validation thereof have been described above. There is one key constraint: the raw VSG input must be high-speed (16-kHz) video because both the HuBERT-Base and HuBERT-Large models [16] were pretrained on speech data encoded as 16-kHz mono WAV files. The PASs must match this sampling rate to ensure appropriate alignment with the

expected input format. Otherwise, feature extraction errors that significantly compromise model performance are to be expected.

However, unlike audio recording, high-speed video recording is resource-intensive. Memory consumption is high, as is the computational cost. This section investigates whether VSG might accept lower-frame-rate videos. Audio waveform upsampling techniques were leveraged when seeking to build a lightweight VSG version.

The experiments investigated this issue using four different upsampling methods. One was a traditional signal processing technique (Bandlimited Sinc Interpolation (BSI) of the widely adopted PyTorch-based Torchaudio library [34]). Additionally, three deep learning approaches were assessed: a deep neural network (DNN) [35], attention-based feature-wise linear modulation (AFiLM) [36], and Phase-Net [37]. A DNN inherently predicts missing high-frequency components in bandwidth-expanded speech and, therefore, enhances model performance during low-frame-rate speech-based ASR tasks. A DNN also addresses the discontinuity between a narrowband input and a reconstructed high-frequency spectrum. AFiLM uses a self-attention mechanism to model long-range temporal dependencies. It is aimed at increasing the fidelity of waveform upsampling. Both of these methods operate directly on 1D PAS waveforms. However, Phase-net uses phase data as the primary inputs. Phase-net was originally designed for video frame interpolation, leveraging phase and amplitude features from complex-valued sub-bands obtained by decomposing each frame using a multi-scale, multi-orientation filter bank (Section 2.2 above). Phase-Net can integrate the obtained phase interpolations with PAS computations, thereby enabling effective PAS upsampling.

BSI upsampling used the parameters of the Torchaudio audio-processing guidelines [34]. These leverage GPU-based acceleration to optimize interpolation, effectively balancing computational efficiency with spectral fidelity. The structural configurations of the deep learning models were those of the original publications [35–37]. Hyperparameter tuning and the training details can be found in the respective references. All four upsampling methods were used to modify original high-speed videos sampled at 8, 4, and 2 kHz, generating 16-kHz PASs via 2 $\times$ , 4 $\times$ , and 8 $\times$  interpolation respectively. The upsampled PASs were used to construct the PASs dataset and then for stage 2 training and evaluation of VSG-Transformer-Large. The CER results are summarized in Table 6.

**Table 6.** The CER results for a VSG using different upsampling techniques and ratios.

Upsampling Methods	2 $\times$ (original 8 kHz)		4 $\times$ (original 4 kHz)		8 $\times$ (original 2 kHz)	
	Development (%)	Test (%)	Development (%)	Test (%)	Development (%)	Test (%)
BSI [35]	22.3	27.1	34.8	41.5	-	-
DNN [36]	14.4	14.8	16.8	17.3	24.7	25.6
AFiLM [37]	14.3	14.8	16.2	16.8	22.5	24.2
Phase-net [38]	13.2	13.6	14.1	15.3	18.4	19.6
-	12.1/12.5 (original 16 kHz)					

As shown in Table 6, at all upsampling ratios, Phase-Net consistently yielded the best VSG performance, with all CERs below 20%. Remarkably, Phase-Net achieved test CERs of 13.6 and 15.3% under 2 $\times$  (original 8 kHz) and 4 $\times$  (original 4 kHz) upsampling. Performance degradation was slight compared to direct training on native 16 kHz data. Performance did not collapse.

Both DNN and AFiLM performed very similarly in terms of the CERs, particularly when recognizing speech at original 4-kHz or higher sampling rates. In contrast, BSI did not perform well. After 8 $\times$  upsampling, VSG-Transformer did not converge during stage 2 training. This may be because, during sampling at 2 kHz, the Nyquist limit restrains the signal bandwidth to below 1 kHz. Such a restricted bandwidth typically captures only the fundamental frequency ( $F_0$ ) and part of the first formants ( $F_1$ ) of vowels and voiced consonants, entirely missing the higher formants such as  $F_2$

and  $F_3$ . Moreover, high-frequency signals from unvoiced consonants such as /f/, /s/, and /k/, which are often above 2 kHz, are completely absent in such low-bandwidth signals [11], significantly impairing recognition accuracy.

Historically, all three deep learning-based upsampling methods sought to incorporate speech features embedded in a high-frequency spectrum into the narrowband speech signal. During training of Phase-Net, DNN, and AFiLM, 16-kHz signals served as the targets of training. In contrast, BSI—a pure numerical interpolation method—fails to recover lost high-frequency features during upsampling. This fundamental distinction largely explains the performance gap between BSI and the deep learning methods. Effective VSG extraction from low-frame-rate data seems to require guidance from high-frame-rate signals during upsampling. However, deep networks featuring learnable upsampling parameters serve as promising directions toward lightweight VSG implementation.

## 5. Conclusions and Future Work

This paper shows that subtle vibrations of everyday objects, caused by ambient sound, can be extracted from video recordings and used to generate text, effectively transforming the objects into vision-based subtitle generators. The VSG pipeline initially acquires pixel-level PASs from the surfaces of objects by PME. These are scalable, abundant robust signals that effectively represent the underlying acoustic features. A pretrained HuBERT-based generative model, termed the VSG-Transformer, is then employed for the PAS-to-text task. The architecture effectively integrates three key components: a convolutional shrinkage module for noise suppression, a pretrained HuBERT-based encoder for extraction of high-level acoustic representations, and an attention-based decoder for autoregressive text generation. Training employs multi-stage transfer learning. The model leverages the knowledge embedded in the pretrained HuBERT and effectively adapts this to recognition of the PAS dataset.

Experimentally, the Base and Large variants of VSG-Transformer successfully accomplished VSG tasks. The CERs of the models were 13.7 and 12.5% respectively. Both models also performed excellently on the conventional ASR task trained on AISHELL-1 with CERs of 6.4 and 6.1% respectively. Comprehensive ablation studies quantitatively evaluated the individual contributions of each module and the overall effectiveness of multi-stage training. Finally, the possibility of lightweight VSG implementation was explored. Low-sampling-rate PASs were upsampled. Four upsampling methods—BSI, DNN, AFiLM, and Phase-Net—that covered both traditional signal processing and deep learning approaches, were tested. Phase-Net upsampling consistently yielded the best performance. Specifically, the VSG-Transformer-Large CERs were 13.6, 15.3, and 19.6% after training on and evaluation of PAS datasets upsampled from original 8- (2×), 4- (4×), and 2-kHz (8×) videos respectively.

In conclusion, VSG fully leverages the flexibility of deep networks and the strength of transfer learning, effectively bridging computer vision-based measurement and natural language processing. The VSG represents a promising new research direction with considerable potential. Future work should explore signal combinations from multiple objects and the integration of inputs from multiple synchronized cameras. Such temporally aligned signals might contain implicit correlations across different sources, further enhancing VSG performance. The rapid advances in smartphone electronic components—particularly the significant increase in camera frame rates (approaching 2 kHz in certain devices)—suggest that VSG deployment in such smartphones is feasible.

**Supplementary Materials:** The following supporting information can be downloaded at the website of this paper posted on Preprints.org.

**Author Contributions:** Conceptualization, Y.W. and X.Z.; methodology, Y.W.; formal analysis, X.Z.; investigation, Y.W.; writing—original draft preparation, Y.W. and X.D.; writing—review and editing, X.D. All authors have read and agreed to the published version of the manuscript.

**Funding:** This work was supported by the National Natural Science Foundation of China (Grant Nos. U2141217).

**Institutional Review Board Statement:** Not applicable.

**Informed Consent Statement:** Not applicable.

**Data Availability Statement:** The data that support the findings of this study are available upon reasonable request from the authors.

**Conflicts of Interest:** The authors declare no conflicts of interest.

## References

1. Davis, A.; Rubinstein, M.; Wadhwa, N.; Mysore, G.; Durand, F.; Freeman, W. The visual microphone: passive recovery of sound from video. *ACM Trans. Graph.* **2014**, *33*, 1–10. <https://doi.org/10.1145/2601097.2601119>.
2. Nassi, B.; Pirutin, Y.; Swissa, R.; Shamir, A.; Elovici, Y.; Zadov, B. Lamphone: Real-time passive sound recovery from light bulb vibrations. *Cryptol. ePrint Arch.* **2020**, *2020*, 4401–4417. <https://eprint.iacr.org/2020/708>.
3. Rothberg, S.; Baker, J.; Halliwell, N. Laser vibrometry: Pseudo-vibrations. *J. Sound Vib.* **1989**, *135*, 516–522. [https://doi.org/10.1016/0022-460X\(89\)90705-0](https://doi.org/10.1016/0022-460X(89)90705-0).
4. Nassi, B.; Swissa, R.; Shams, J.; Zadov, B.; Elovici, Y. The little seal bug: Optical sound recovery from lightweight reflective objects. In Proceedings of the IEEE Security and Privacy Workshops. **2023**, 298–310. <https://doi.org/10.1109/SPW59333.2023.00032>.
5. Nassi, B.; Pirutin, Y.; Galor, T.; Elovici, Y.; Zadov, B. Glowworm attack: Optical tempest sound recovery via a device's power indicator led. In Proceedings of the ACM SIGSAC Conference on Computer and Communications Security. **2021**, 1900–1914. <https://doi.org/10.1145/3460120.3484775>.
6. Kwong, A.; Xu, W.; Fu, K. Hard drive of hearing: Disks that eavesdrop with a synthesized microphone. In Proceedings of the IEEE Symposium on Security and Privacy. **2019**, 905–919. <https://doi.org/10.1109/SP.2019.00008>.
7. Zhang, D.; Guo, J.; Jin, Y.; Zhu, C. Efficient subtle motion detection from high-speed video for sound recovery and vibration analysis using singular value decomposition-based approach. *Opt. Eng.* **2017**, *56*, 094105. <https://doi.org/10.1117/1.OE.56.9.094105>.
8. Zhang, D.; Guo, J.; Lei, X.; Zhu, C. Note: sound recovery from video using svd-based information extraction. *Rev. Sci. Instrum.* **2016**, *87*, 086111. <https://doi.org/10.1063/1.4961979>.
9. Guri, M.; Solewicz, Y.; Daidakulov, A.; and Elovici, Y. SPEAKE(a) R: Turn speakers to microphones for fun and profit. In Proceedings of the 11th USENIX Workshop on Offensive Technologies. **2017**. <https://doi.org/10.48550/arXiv.1611.07350>.
10. Michalevsky, Y.; Boneh, D.; Nakibly, G. Gyrophone: Recognizing speech from gyroscope signals. In Proceedings of the 23rd USENIX Security Symposium. **2014**, 1053–1067.
11. Long, Y.; Naghavi, P.; Kojusner, B.; Butler, K.; Rampazzi, S.; Fu, K. Side eye: Characterizing the limits of pov acoustic eavesdropping from smartphone cameras with rolling shutters and movable lenses. In Proceedings of the IEEE Symposium on Security and Privacy. **2023**, 1857–1874. <https://doi.org/10.1109/SP46215.2023.10179313>.
12. Zhang, L.; Pathak, P.; Wu, M.; Zhao, Y.; and Mohapatra, P. AccelWord: Energy efficient hotword detection through accelerometer. In Proceedings of the 13th Annual International Conference on Mobile Systems, Applications, and Services. **2015**, 301–315. <https://doi.org/10.1145/2742647.2742658>.
13. Han, J.; Chung, A.J.; Tague, P. PitchIn: Eavesdropping via intelligible speech reconstruction using non-acoustic sensor fusion. In Proceedings of the 16th ACM/IEEE International Conference on Information Processing in Sensor Networks. **2017**, 181–192. <https://doi.org/10.1145/3055031.305508>.
14. Wang, G.; Zou, Y.; Zhou, Z.; Wu, K.; Ni, L. We can hear you with Wi-Fi! *IEEE Trans. Mobile Comput.* **2016**, *15*, 2907–2920. <https://doi.org/10.1109/TMC.2016.2517630>.
15. Devlin, J.; Chang, M.; Lee, K.; Toutanova, K. BERT: Pre-training of deep bidirectional transformers for language understanding. In Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies. **2019**, 4171–4186. <https://doi.org/10.18653/v1/N19-1423>.

16. Hsu, W.; Bolte, B.; Tsai, Y.; Lakhotia, K.; Salakhutdinov, R.; Mohamed, A. Hubert: Self-supervised speech representation learning by masked prediction of hidden units. *IEEE/ACM Trans. Audio, Speech, Lang. Process.* **2021**, *29*, 3451–3460. <https://doi.org/10.1109/TASLP.2021.3122291>.
17. Yang, Z.; Dai, Z.; Yang, Y.; Carbonell J.; Salakhutdinov, R.; Le, Q. XLNet: Generalized autoregressive pretraining for language understanding. In Proceedings of the Neural Information Processing Systems. **2019**, 5754–5764.
18. Liu, Y., Ott, M.; Goyal. N.; Du, J.; Joshi, M.; Chen, D.; Levy, O.; Lewis, M.; Zettlemoyer, L.; Stoyanov, V. Roberta: A robustly optimized bert pretraining approach. **2019**. arXiv:1907.11692.
19. Bu, H.; Du, J.; Na, X.; Wu, B.; Zheng, H. AISHELL-1: An open-source Mandarin speech corpus and a speech recognition baseline. In Proceedings of the 20th Conference of the Oriental Chapter of the International Coordinating Committee on Speech Databases and Speech I/O Systems and Assessment. **2017**, 1–5. <https://doi.org/10.1109/ICSDA.2017.8384449>.
20. Fleet, D.; Jepson A. Computation of component image velocity from local phase information. *Int. J. Comput. Vis.* **1990**, *5*, 77–104. <https://doi.org/10.1007/BF00056772>.
21. Gautama T.; VanHulle, M. A phase-based approach to the estimation of the optical flow field using spatial filtering. *IEEE Trans. Neural Netw.*, **2002**, *13*, 1127–1136. <https://doi.org/10.1109/TNN.2002.1031944>.
22. Freeman, W.; Adelson, E. The design and use of steerable filters. *IEEE Trans. Pattern Anal. Mach. Intell.* **1991**, *13*, 891–906. <https://doi.org/10.1109/34.93808>.
23. Chou, J.; Chang, C.; Spencer J. Out-of-plane modal property extraction based on multi-level image pyramid reconstruction using stereophotogrammetry. *Mech. Syst. Signal Process.* **2022**, *169*, 108786. <https://doi.org/10.1016/j.ymsp.2021.108786>.
24. Wadhwa, N.; Rubinstein, M.; Durand, F.; Freeman, W. Phase based video motion processing. *ACM Trans. Graph.* **2013**, *32*, 1–10. <https://doi.org/10.1145/2461912.2461966>.
25. Isogawa, K.; Ida, T.; Shiodera, T.; Takeguchi, T. Deep shrinkage convolutional neural network for adaptive noise reduction. *IEEE Signal Process. Lett.* **2018**, *25*, 224–228. <https://doi.org/10.1109/LSP.2017.2782270>.
26. Zhao, M.; Zhong, S.; Fu, X.; Tang, B.; Pecht, M. Deep residual shrinkage networks for fault diagnosis. *IEEE Trans. Ind. Informat.* **2020**, *16*, 4681–4690. <https://doi.org/10.1109/TII.2019.2943898>.
27. Zhang, B.; Lv, H.; Guo P.; Shao Q.; Yang C.; Xie L. Wenetspeech: A 10000 hours multi-domain mandarin corpus for speech recognition. In Proceedings of the IEEE International Conference on Acoustics, Speech and Signal Processing. **2022**, 6182–6186. <https://doi.org/10.1109/ICASSP43922.2022.9746682>.
28. Joshi, M.; Chen, D.; Liu, Y.; Weld, D.; Zettlemoyer, L.; Levy, O. SpanBERT: Improving pre-training by representing and predicting spans. *Trans. Assoc. Comput. Linguist.* **2020**, *8*, 64–77. [https://doi.org/10.1162/tacl\\_a\\_00300](https://doi.org/10.1162/tacl_a_00300).
29. Baeovski, A.; Zhou, H.; Mohamed, A.; Auli, M.; wav2vec 2.0: A framework for self-supervised learning of speech representations. **2020**. arXiv:2006.11477.
30. Paszke, A.; Gross, S.; Massa, F.; Lerer, A.; Bradbury, J.; Chanan, G.; Killeen, T.; Lin, Z.; Gimelshein, N.; Antiga, L.; et al. PyTorch: An Imperative Style, High-Performance Deep Learning Library. *Adv. Neural Inf. Process. Syst.* **2019**.
31. Kingma, D.; Ba, J. Adam: A method for stochastic optimization. **2015**, arXiv:1412.6980.
32. Chorowski J.; Jaitly, N. Towards better decoding and language model integration in sequence to sequence models. In Proceedings of the Interspeech. **2017**, 523–527. <https://doi.org/10.48550/arXiv.1612.02695>.
33. Wu, Y.; Schuster, M.; Chen, Z.; Le, Q.; Norouzi, M.; Macherey, W.; Krikun, M.; Cao, Y.; Gao, Q.; Macherey K.; et al. Google’s neural machine translation system: Bridging the gap between human and machine translation. **2016**, arXiv:1609.08144.
34. Yang, Y.; Hira, M.; Ni, Z.; Astafurov, A.; Chen, C.; Puhersch, C. Torchaudio: Building blocks for audio and speech processing. In Proceedings of the IEEE International Conference on Acoustics, Speech and Signal Processing. **2022**, 6982–6986. <https://doi.org/10.1109/ICASSP43922.2022.9747236>.
35. Li, K.; Huang, Z.; Xu, Y.; Lee C. DNN-based speech bandwidth expansion and its application to adding high-frequency missing features for automatic speech recognition of narrowband speech. In Proceedings of the Interspeech. **2015**, 2578–2582.

36. Rakotonirina, N. Self-attention for audio super-resolution. In Proceedings of the IEEE International Workshop on Machine Learning for Signal Processing, 2021, 1–6. <https://doi.org/10.1109/MLSP52302.2021.9596082>.
37. Meyer, S.; Djelouah, A.; McWilliams, B.; Sorkine-Hornung, A.; Gross, M.; Schroers, C. Phasenet for video frame interpolation. In Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. 2018, 498–507.

**Disclaimer/Publisher's Note:** The statements, opinions and data contained in all publications are solely those of the individual author(s) and contributor(s) and not of MDPI and/or the editor(s). MDPI and/or the editor(s) disclaim responsibility for any injury to people or property resulting from any ideas, methods, instructions or products referred to in the content.