
SILENCE: Protecting privacy in offloaded speech understanding on wimpy devices

Abstract

1 Speech serves as a ubiquitous input interface for embedded mobile de-
2 vices. Cloud-based solutions, while offering powerful speech understanding
3 services, raise significant concerns regarding user privacy. To address this,
4 disentanglement-based encoders have been proposed to remove sensitive infor-
5 mation from speech signals without compromising the speech understanding func-
6 tionality. However, these encoders demand high memory usage and computation
7 complexity, making them impractical for resource-constrained wimpy devices.
8 Our solution is based on a key observation that speech understanding hinges on
9 long-term dependency knowledge of the entire utterance, in contrast to privacy-
10 sensitive elements that are short-term dependent. Exploiting this observation, we
11 propose SILENCE, a lightweight system that selectively obscuring short-term de-
12 tails, without damaging the long-term dependent speech understanding perfor-
13 mance. The crucial part of SILENCE is a differential mask generator derived from
14 interpretable learning to automatically configure the masking process. We have
15 implemented SILENCE on the STM32H7 microcontroller and evaluate its efficacy
16 under different attacking scenarios. Our results demonstrate that SILENCE offers
17 speech understanding performance and privacy protection capacity comparable to
18 existing encoders, while achieving up to $53.3\times$ speedup and $134.1\times$ reduction in
19 memory footprint.

20 1 Introduction

21 **Privacy concern for cloud speech service** The volume of speech data uploaded to the cloud for
22 spoken language understanding (SLU) is steadily increasing [1, 12, 2], particularly in ubiquitous
23 wimpy devices where textual input is inconvenient [41, 17, 3], e.g., home automation devices [32],
24 smartwatches [37], telehealth sensors [22] and smart factory sensors [29]. However, exposing raw
25 speech signal to the cloud raises privacy concerns [42]. It was revealed that contractors regularly
26 listened to confidential details in Siri recordings to improve its accuracy [4]. This included private
27 discussions, medical information, and even intimate moments.

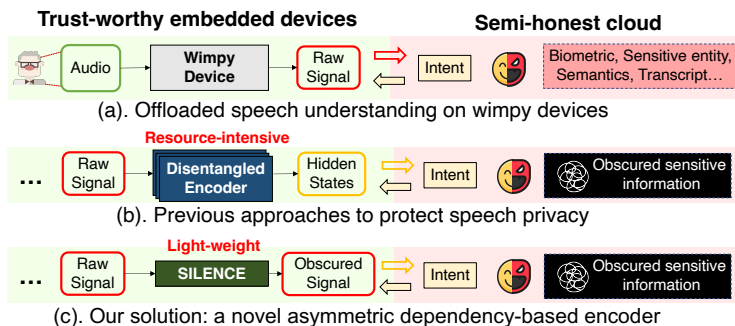


Figure 1: Illustration of offloaded speech understanding on wimpy devices and its privacy protection.

28 There are many aspects of potential privacy leakage in cloud-based SLU. Among them: biometric
29 or contextual privacy leakage have been well studied and somewhat solved by removing information
30 relevant to such tasks without compromising the SLU accuracy [18, 35]; transcript protection (espe-
31 cially sensitive entities) is more challenging since it is deeply entangled with the SLU task itself. As
32 shown in Figure 1, this paper focus on ensuring that cloud-based systems could efficiently classify
33 the intent of SLU task (e.g., scheduling appointments or controlling home devices) while refraining
34 from identifying the concrete entities (e.g., unintended names or passwords) in the spoken utterance,
35 i.e., high word error rate (WER) of Automatic Speech Recognition (ASR) task. This is also a setting
36 commonly used in speech privacy protection [44, 10, 16, 42, 15].

37 **Prior approaches** A prevalent method for private speech processing is employing *encoders*¹ based
38 on disentanglement representation learning [44, 10, 28, 34], as illustrated in Figure 1(b). Those en-
39 coders extract the speech representations using pre-trained acoustic models, e.g., wav2vec [40, 10],
40 conformer [26, 34] and Preformer [20, 44]. Furthermore, they promote representation disentangle-
41 ment through adversarial training [25]. For example, PPSLU [44] uses a 12-layer transformer-based
42 Preformer as its encoder.

43 As a result, disentanglement-based encoders still demand considerable computational resources,
44 often exceeding tens of GFLOPs, to achieve effective disentanglement [11]. They are also memory-
45 intensive, often comprising tens of millions of parameters. Consequently, they are unsuitable for
46 embedded devices with limited memory. Moreover, it takes time-consuming adversarial training to
47 disentangle the encoded representation for each specific SLU task. This aspect limits the flexibility
48 and scalability for emerging SLU tasks. More motivating details will be presented in §2.2.

49 In this paper, we aim to achieve the real-time, privacy-preserving offloading of speech understanding
50 task on wimpy devices like STM32H7 microcontroller [5] with only 1MB RAM. This goal neces-
51 sitates a novel encoder design that must be both lightweight and effective in filtering out sensitive
52 information, as illustrated in Figure 1(c).

53 **Our solution** We therefore present SILENCE, a **SImpLe ENCOder** designed for efficient privacy-
54 preserving SLU offloading. It is based on the *asymmetric dependency* observation: SLU intent
55 extraction (e.g., scenario identification) typically requires only long-term dependency knowledge
56 across the entire utterance, while ASR task (e.g., recognizing individual words or phrases) needs
57 short-term dependency, as confirmed by our experiments in §3.1. Based on it, SILENCE strategically
58 partitions the utterance into several segments, selectively masking out the majority to enhance pri-
59 vacy by obscuring short-term details, without significantly damaging the long-term dependencies.
60 The processed audio waveform is then transmitted to the cloud for SLU intent analysis. Addition-
61 ally, we integrate a differential mask generator, inspired by interpretable learning methods [19], to
62 optimize performance by automatically identifying how many and which segments to mask.

63 **Results** We deploy SILENCE on the STM32H7 microcontroller [5] and assess its performance
64 using the SLURP dataset [13] in both black-box and white-box attack environments. SILENCE
65 achieves 81.2% intent classification accuracy on SLURP, surpassing previous privacy-preserving
66 SLU systems by up to 8.3%. Regarding privacy protection, SILENCE offers comparable security
67 to earlier systems, with a word error rate of up to 81.6% and an entity error rate of 90.7% under
68 malicious ASR attacks. Even against white-box attacks, where attackers are strongly assumed to
69 have the same encoder structure and weights as SILENCE, plus partial data from malicious clients,
70 SILENCE maintains 67.3% word error rate and 64.3% entity error rate. Additionally, SILENCE
71 proves to be resource-efficient and feasible for wimpy devices, using only 394.9KB of memory
72 and taking just 912.0ms to encode a 4-second speech signal. Integrated with RPI-4B for a fair
73 comparison, SILENCE uses up to $134.1\times$ less memory and operates up to $53.3\times$ faster than prior
74 systems. The accuracy of SILENCE is only 7% lower than unprotected SLU systems.

75 **Contribution** We have made the following contributions.

- 76 • Based on the observation of asymmetric dependency between SLU and ASR tasks, we
77 propose SILENCE, a simple yet effective encoder system for privacy-preserving SLU of-
78 floading.

¹Note that these encoders are not specifically transformer encoders; rather, they can be implemented using any NNs to encode speech signals.

- 79 • We are the first to retrofit interpretable learning methods to automatically configure the
80 masking process for a better balance between privacy and utility in speech understanding
81 tasks.
- 82 • We evaluate SILENCE on a wimpy microcontroller unit and demonstrate its effectiveness
83 under various attack scenarios.

84 2 Related Work and Background

85 2.1 Privacy-preserving SLU

86 Spoken Language Understanding (SLU) is a critical component of modern voice-activated systems,
87 responsible for interpreting human speech and translating it into structured, actionable commands.
88 For instance, when a user says, "Set a meeting for tomorrow at 10 AM," the SLU system might map
89 this to a structured intent such as `{scenario: Calendar, action: Create_entry}`.

90 **Evolution of SLU Systems** The evolution of SLU systems has seen a shift from traditional two-
91 component systems, comprising ASR and Natural Language Understanding (NLU), to modern end-
92 to-end neural networks [39, 27]. These advanced systems bypass the intermediate textual represen-
93 tation and directly map speech signals to their semantic meaning, enhancing efficiency and reducing
94 error propagation. A typical end-to-end SLU model features an encoder, often with convolution and
95 attention-based elements, and a decoder, including a transformer decoder and a connectionist tem-
96 poral classification decoder. Many SLU systems incorporate encoders from pre-trained ASR models
97 like HuBERT [45], replacing the original ASR decoder with one tailored for SLU tasks.

98 **Threat Model** Our threat model aligns with prior work [44, 10] where users (the victims) actively
99 offloads their audio data to the cloud server (the adversary) for intended SLU tasks. Upon receiving
100 the data, the adversary may employ automatic speech recognition to transcribe the audio and identify
101 private entities [16, 42, 15]. Note that the transcriptions are often exceedingly detailed, containing
102 much more information than the users intend to disclose. The goal of this paper is to ensure that
103 the victims can reliably obtain the predefined SLU intent from the adversary, while preserving the
104 adversary from discerning sensitive details or private entities in the transcript.

105 For instance, home pods might capture recordings of confidential daily interactions alongside ex-
106 plicit commands, presenting a paradigmatic case for SILENCE. Without SILENCE, over 80% of our
107 private daily conversations could be automatically recognized and stored for unforeseen usage as
108 will be analyzed in §5.1.

109 2.2 Inefficiency of Existing Approaches

110 **Privacy-preserving methods** Crypto-based approaches, such as HE [48] and MPC [24], have been
111 proposed to provide encrypted computation. Unfortunately, they are technically slow and thus im-
112 practical for deployment on wimpy audio devices due to the significant increase in computation
113 and communication complexity. For example, MPC-based PUMA [21] takes 5 minutes to com-
114 plete one token inference, which is far too slow for real-time. Voice conversion is another method
115 to protect speech content. `PREECH` [9] integrates voice conversion with GPT-based generated noise
116 to protect privacy, but it is far from feasible for deployment on wimpy devices. Traditional periph-
117 eral devices, such as ultrasonic microphone jammers (UMJ), are designed to obscure raw speech by
118 inserting non-linearity noise, thereby preventing illegal eavesdropping[23, 15]; however, they also
119 corrupt speech semantics as well. A emerging and prevailing strategy is disentangling-based en-
120 coders [10, 44, 28]; they aim to create a disentangled and hierarchical representation of the speech
121 signal devoid of sensitive data. But we reveal their performance issue next.

122 We conduct preliminary experiments to measure the resource consumption of the disentangling-
123 based encoder of a pre-trained SLU model on a Raspberry Pi 4B (RPI-4B) [6] and Jetson TX2
124 (TX2) [7]. Our key observation is that disentangling-based privacy-preserving SLU system is too
125 resource-intensive for practical deployment. As illustrated in Figure 2, a disentanglement encoder
126 consumes 648.7MB memory and 12.8s for complete one inference on RPI-4B. Even in the strong
127 TX2 with GPU, the encoder still takes 593.0ms to complete one inference. Considering the network
128 latency, the end-to-end latency of the disentangling-based SLU offloading system only saves 0.7%

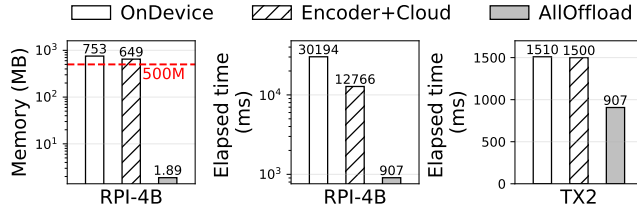


Figure 2: Cost of disentangling-based encoders [44] for a 4-second audio inference.

129 wall-clock time compared to the OnDevice inference without offloading, with a similar memory
 130 footprint over 500M.

131 **Implications** Disentangling-based encoders is slow and memory-intensive due to the complex en-
 132 coder structure designed to separate sensitive information from the speech signal. Given the limited
 133 resource of wimpy devices, it is not practical for common privacy-preserving SLU scenarios. To
 134 enable practical privacy-preserving SLU, the encoder structure and the inference process need to be
 135 simplified.

136 3 SILENCE Design

137 3.1 System Design and Rationales

138 We introduce SILENCE to efficiently scrub raw audio for privacy-preserving SLU, as depicted in
 139 Figure 3. The key idea of SILENCE is simple and novel: it masks out a portion of audio segments
 140 before sending them to the cloud for SLU tasks. This design is based on an unique observation
 141 shown in Figure 4(c): when a portion of audio segments is masked out, the ASR model becomes
 142 incapable to recognize the phonemes in the masked frames, while the SLU model can still recognize
 the intent.

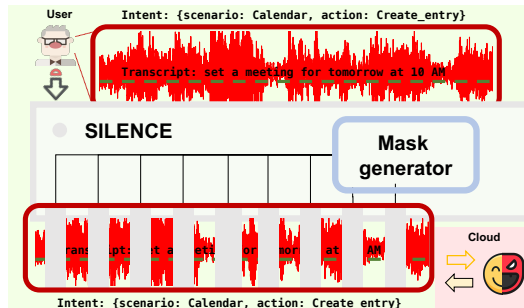


Figure 3: SILENCE overview. Red hard line represents the long-term dependency, while the green
 143 dotted line represents the short-term dependency.

144 **Design rationale** Why is SILENCE able to protect the sensitive entity privacy while maintaining
 145 SLU accuracy? This capability is rooted in the *asymmetrical dependency* between the ASR and
 146 SLU task.

147 Speech is composed of many meta phonemes, and the generation of a single meta phoneme depends
 148 on its adjacent frame [42]. *Dependency* is defined as the length of frame that a model's output
 149 depends on. Figure 4(a) shows each phoneme is mainly dependent on a few frames, indicating short-
 150 term dependency. This phenomenon is referred to as "peaky behavior" in the ASR literature [47]. In
 151 contrast, an SLU model utilizes an attention-based decoder [45] to capture the relationship between
 152 the entire utterance and the intent, implying that the intent is long-term dependent on the whole
 153 utterance.

154 Formally, SILENCE is a simple encoder based on asymmetrical dependency-based masking. This
 155 simple masking encoder is defined as: $\hat{x} = x \odot \mathbb{Z}$, where x is the input audio signal, \odot represents
 156 the element-wise multiplication, \hat{x} is the masked audio signal and \mathbb{Z} is the binary masking vector
 157 with the same dimension as x . \mathbb{Z} consists of k uniform portion, with all 0s or 1s in one portion

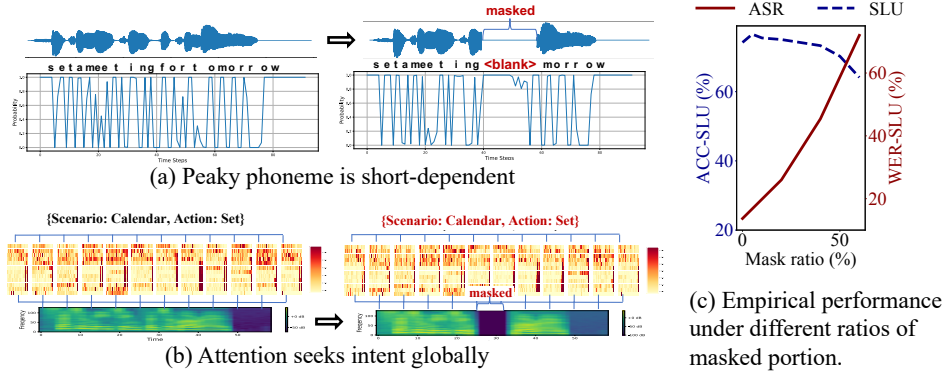


Figure 4: Foundation of SILENCE: asymmetrical dependency. (a). ASR task is short-term dependent on the peaky phoneme probability. (b). SLU task is long-term dependent on knowledge from the whole utterance. (c). Empirical results.

158 to mask-out or preserve the complete adjacent frames, respectively. This simple encoder forms
 159 the basis of SILENCE’s efficiency and privacy-preservation capacity, enabling secure offloading of
 160 speech understanding tasks on wimpy devices.

161 **The configuration challenges:** Figure 4(c) demonstrates that the ratio of masked portion plays
 162 a crucial role in balancing the privacy (WER-ASR) and utility (ACC-SLU). Currently, SILENCE
 163 employs a trivial masking mechanism, necessitating clients to undertake a time-intensive hyper-
 164 parameter adjustment about the extent and location of masking. Incorrect masking configurations
 165 can result in significant loss of global long-term dependency, negatively affecting SLU accuracy,
 166 or insufficient masking of sensitive information, thus compromising privacy. Therefore, we face
 167 critical questions: how many and which portions should be masked?

168 3.2 Online Configurator for SILENCE

169 To address these challenges, we derive a differential mask generator from the interpretable learn-
 170 ing [19] as an online configurator for SILENCE. This automatically generate the masking vector \mathbb{Z} .
 171 The mask generator is trained to identify how many and which portions to mask, optimizing the
 172 privacy-utility balance.

173 **Differentiable mask generator** The configurator model aims to minimize the discrepancy between
 174 masked and original output by generating a mask \mathbb{Z} . Formally, we define the number of unmasked
 175 portions as \mathcal{L}_0 loss:

$$\mathcal{L}_0(\phi, x) = \sum_{i=1}^n \mathbf{1}_{[\mathbb{R}_{\neq 0}]}(\mathbb{Z}_i) \quad (1)$$

176 where ϕ is the mask generator, $\mathbf{1}(\cdot)$ is the indicator function. We minimize \mathcal{L}_0 for dataset \mathcal{D} , ensuring
 177 that predictions from masked inputs resemble those from the origin model:

$$\min_{\phi} \sum_{x \in \mathcal{D}} \mathcal{L}_0(\phi, x) \quad (2)$$

$$\text{s.t. } D_{\star}[y||\hat{y}] \leq \gamma \quad \forall x \in \mathcal{D} \quad (3)$$

178 where $\hat{y} = f(\hat{x})$, y is the tokenized label, $D_{\star}[y||\hat{y}]$ is the KL divergence and the margin $\gamma \in \mathbb{R}_{>0}$ is
 179 a hyperparameter.

180 Given that \mathcal{L}_0 is discontinuous and has zero derivative almost everywhere, and the mask generator ϕ
 181 requires a discontinuous output activation (like a step function) for binary masks, we utilize a sparse
 182 relaxation to binary variables [30, 14] instead of the binary mask during training.

183 **Holistic workflow** As shown in Figure 5, SILENCE encompasses two phases:

184 (1) *Offline phase:* **(1a)** First, SILENCE trains a differentiable mask generator. The client selects a
 185 mask generator model, potentially a submodule of a pre-trained ASR model, such as HuBERT’s

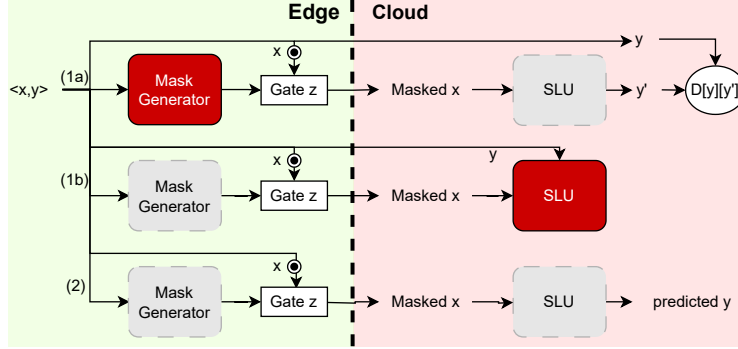


Figure 5: SILENCE workflow. (1) *Offline phase*: (1a) Training mask generator and (1b) adapting cloud SLU model to it; (2) *Online phase*: Conducting cloud inference with the masked x . Only masked input audio x and insensitive intent label y are exposed to the cloud.

186 CNN feature extractor. A small gate model is then integrated with this submodule. The combined
 187 model processes the input audio and generates a mask. This mask selectively conceals parts of the
 188 input, ensuring retention of only vital SLU information while hiding sensitive data. The masked
 189 input is then forwarded to either a trusted cloud service or a local SLU model for obtaining masked
 190 output. The mask generator is fine-tuned to minimize the discrepancy between the masked output
 191 logits and the original intent, as defined in Equation (1-3).

192 (1b) Second, SILENCE adapts the cloud model. Here, the client forwards the masked input and a
 193 specific SLU intent (e.g., "set alarm") to the cloud-based SLU model. The model undergoes fine-
 194 tuning to adapt to the masked inputs. This process includes adjusting the model parameters for
 195 accurate recognition and response to SLU commands based on the masked input.

196 (2) *Online phase*: In online speech understanding, the client sends the masked input to the cloud
 197 SLU model. Using the adapted model, the cloud-based SLU accurately identifies and executes the
 198 intended SLU action or response.

199 **Configurator cost analysis** Training the differentiable mask generator is affordable for the client.
 200 Our experiments indicate that convergence is achieved with approximately 200 audio samples,
 201 equivalent to 600 seconds of audio. This process takes up to 30 seconds on an A40 GPU. Adapting
 202 the SLU model to each mask generator is a one-pass effort. This adaptation is relatively trivial, espe-
 203 cially when starting from a fine-tuned SLU model rather than building from scratch. This aspect of
 204 the process incurs minimal cost compared to the training of the cloud SLU model. Moreover, these
 205 costs can be amortized over a large number of edge users in the long run, making it an economically
 206 viable solution.

207 **Remark** Note that the mask generator is not developed for tagging sequences at a semantic level.
 208 Rather, its design focuses on identifying segments that are more relevant to the SLU task. This task is
 209 essentially a relatively straightforward binary classification problem, which is proven to be effective
 210 in prior interpretable learning literature [19, 14] and light-weight enough for real-time inference.

211 4 Implementation and Methodology

212 We have fully implemented the SILENCE prototype atop SpeechBrain [38], a PyTorch-based and
 213 unified speech toolkit. As prior work [45], we use SpeechBrain to train the differential mask gener-
 214 ator and simulate the cloud training process. After that, we deploy the trained mask generator into
 215 the embedded devices and evaluate the end-to-end performance.

216 **Hardware and environment** Offline training is simulated on a server with 8 NVIDIA A40 GPUs.
 217 The trained mask generator is deployed into the STM32H7 [5] or Raspberry PI 4 (RPI-4B) [6].
 218 STM32H7 is a wimpy microcontroller with 1MB RAM. RPI-4B is a popular development board
 219 with 4GB RAM. We embed the approaches not feasible to fit in the STM32H7 into the RPI-4B.

220 **Models** We design four types of mask generator structures: (1) Random: a random binary vector
 221 generator with 50% portion masked; (2) SILENCE-S: a learnable mask generator with only one MLP

222 gate; (3) SILENCE-M: a learnable mask generator with one HuBERT encoder layer and the gate; (4)
223 SILENCE-L: a learnable mask generator with three HuBERT encoder layers and the gate. As for the
224 cloud SLU model, we simulate it using the SoTA end-to-end SLU model [45]. It replaces the ASR
225 decoder of pre-trained HuBERT with SLU attentional decoder.

226 **Dataset and Metrics** We run our experiments on SLURP [13] with 102 hours of speech. SLURP’s
227 utterances are complex and closer to daily human speech. We select scenario classification accu-
228 racy to measure the SLU understanding performance (ACC-SLU). Following prior work [44], we
229 choose large-scale English reading corpus LibriSpeech [33] for a multi-task protection scenario.
230 In the multi-task protection scenario, not only the SLU command utterance (SLURP) but also the
231 background or the subsequent utterance (LibriSpeech) are uploaded to the cloud. WER is used to
232 measure the attack performance. More specifically, we utilize WER-SLU to measure the attacker’s
233 capacity to recognize the word information in the uploaded SLU audio itself, and WER-ASR as
234 the WER of recognized accompanying audio, i.e., LibriSpeech dataset. We also report the private
235 entity recognition error rate (EER) to ensure that the cloud model is not able to recognize the private
236 information in the speech signal. As for latency, we sequentially fed test audios into the local model
237 without any window processing² and recorded the average forward time as the local execution time.

238 **Baselines** We compare SILENCE to the following alternatives: (1) `OnDevice` means the cloud SLU
239 model is downloaded and run locally on the client device. (2) `AllOffload` means the raw audio
240 is uploaded to the cloud for SLU inference. (3) VAE [10] is the vanilla variational auto-encoder
241 method that uses adversarial training to disentangle the private information from speech signal. (4)
242 PPSLU [44] is the state-of-the-art disentangling-based SLU privacy-preserving system, which uses
243 12 transformer layers to separate the SLU information into a part of the hidden layer and only sends
244 those hidden layers to the cloud for SLU inference.

245 **Attack scenarios.** We use three attacks encompassing both black-box and white-box attacks:
246 (1) `Azure` represents a black-box attacker scenario, in which the masked audio is transmitted to
247 Azure [31] for automatic speech recognition. (2) `Whisper` simulates a SoTA cloud-based ASR
248 model. This black-box attacker uses the pre-trained *Whisper.medium.en* model [36], directly
249 downloaded from HuggingFace [46]. (3) `Whisper(White-box)` constitutes a white-box attack.
250 Here, we hypothesize that certain users are malicious and disclose the mask generator’s structure
251 and weights, along with their own audio data, to the `Whisper` attack model. `Whisper(White-box)`
252 then utilizes this collected data from malicious users to adapt the pre-trained *Whisper.medium.en*
253 model to the specific masking pattern.

254 **Hyper-parameters** During the offline phase in Figure 5, we use the Adam optimizer with a learning
255 rate of $1e-5$ and a batch size of 4. For the inference step, we use the batch size of 1 to simulate the
256 real streaming audio input scenario. The end-to-end cloud SLU latency is measured by invoking
257 Azure APIs following previous work [43]. KL threshold λ is set as 0.15 for all mask generators.
258 Attack model is set as `Whisper` without special declaration.

259 5 Evaluation

260 5.1 End-to-end performance

261 **SILENCE achieves comparable accuracy performance and privacy protection capacity to pre-**
262 **vious encoders.** As shown in Figure 6, we compare the accuracy of SILENCE with all baselines.
263 `OnDevice` offloads no signals to the cloud and thus has the best privacy protection (WER=100).
264 It is observed that SILENCE could achieve up to 81.1% accuracy, with less than 7% accuracy loss
265 compared to unprotected `AllOffload` and local `OnDevice` SLU model. Its rationale is that we
266 mainly mask the short-dependent frames that does not significantly affect the SLU performance.
267 We also compare the performance of SILENCE with the SoTA privacy-preserving SLU system, i.e.,
268 PPSLU [44]. SILENCE achieves 7.2% higher accuracy than PPSLU which tries to apply complex non-
269 linear transformation to the hidden layer to prevent malicious re-construction, but this might also
270 damage part of the SLU information. In terms of privacy preservation, our learnable mask generator
271 achieves up to 78.6% WER using SILENCE-L, indicating a privacy-preserving capacity on par with

²The average duration of test SLU snippets is 2.8 seconds, with a maximum of 21.5 seconds, which is shorter than the maximum input window of speech models (e.g., 30 seconds for `Whisper` [36]).

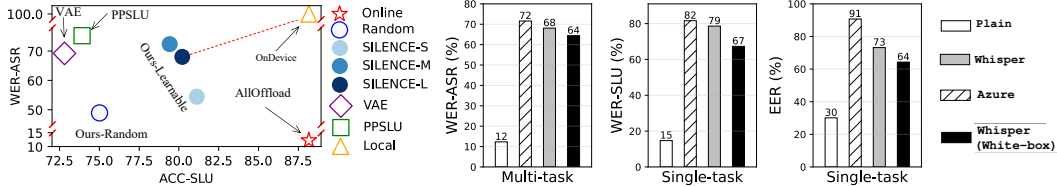


Figure 6: Performance of different privacy-preserving SLU approaches.

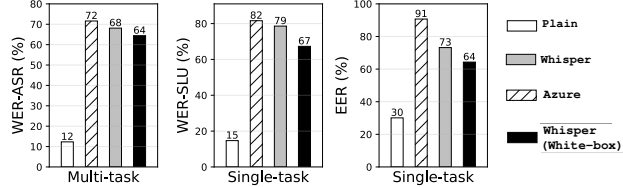


Figure 7: SILENCE privacy-preserving capacity under different attack models.

272 PPSLU. Furthermore, we complete the inference with much lower delays and memory footprint as
 273 will be shown in Figure 9.

274 **SILENCE is resistant to different attack models.** As illustrated in Figure 7, SILENCE increases the
 275 SLU-WER from 14.7% to 78.6% under the attack model Whisper. As for the online attack model
 276 Azure, SILENCE increases the SLU-WER from 14.7% to 81.6%. According to our returned service
 277 details, we find that over 50% of the sent audios are tagged as "ResultReason.NoMatch", which
 278 means audios are recognized as null utterances by the Azure ASR model. Whisper (White-box)
 279 is a white-box attack model, which means the attacker has the same mask generator structure and
 280 weights as the SILENCE. We still achieve more than 50% SLU-WER under this attack model. This
 281 is because even Whisper (White-box) is fine-tuned to fill some of the missing frames, it still could
 282 not recover the private missing frames. Because masking the short-dependent frames fundamentally
 283 destroys the raw audio signal. It is not possible to re-construct the phoneme without knowing any
 284 speech information. In the last subfigure, we show the high entity error rate to demonstrate that the
 285 private entity is not leaked.

286 **SILENCE scales to better privacy-accuracy trade-off with a larger mask generator.** We explore
 287 the impact of the threshold γ of SILENCE under different mask generator structures. As shown in
 288 Figure 8, the threshold γ controls the trade-off between the privacy and utility. When γ is small,
 289 the mask generator is more conservative, leading to higher the utility a lower the masking portion.
 290 As we have discussed in Section 3, a lower rate of masking portions leads to higher possibility of
 291 privacy entity leakage. When γ is large, the mask generator is more aggressive, enhancing privacy.
 292 Another way to achieve more practical privacy-utility balance is using a more complex mask gener-
 293 ator structure, e.g., SILENCE-L. It achieves higher utility with the same privacy level compared to
 294 SILENCE-S, albeit with less efficiency, as shown in § 5.2.

295 5.2 System cost

296 SILENCE protects the private entities efficiently as shown in Figure 9. Different from prior encoders
 297 using complex disentanglement model, SILENCE only requires a light-weight mask generator to
 298 scrub the private information. The size of this generator varies according to different mask gener-
 299 ator structures. For the smallest mask generator, SILENCE-S, it only requires a 394.9KB memory
 300 footprint, and could successfully embed into the wimpy STM32H7 with 2MB RAM. SILENCE is
 301 efficient not only in terms of memory footprint but also in latency. SILENCE-S completes the local
 302 encoding with only 912.2ms on the wimpy STM32H7. For a fair comparison, we embed SILENCE-S

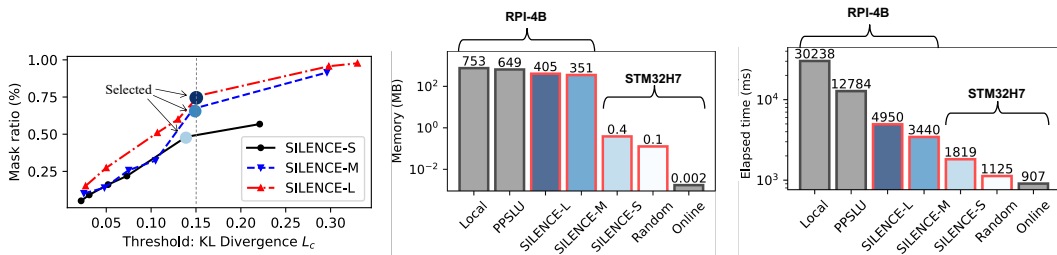


Figure 8: Effect of threshold with different mask generators

(a) Memory footprint

(b) End-to-end latency

Figure 9: Comparison of resource cost in different SLU approaches. Ours are highlighted in red.

303 into RPI-4B and find that it is $18.1\times$ faster and $134.1\times$ less memory footprint than PPSLU. Even with
304 the strong mask generator SILENCE-L, SILENCE achieves up to $7.5\times$ lower encoding latency and
305 consumes $1.9\times$ less memory compared to OnDevice.

306 6 Conclusion and Discussions

307 SILENCE is an efficient and privacy-preserving end-to-end SLU system based on the asymmetrical
308 dependency between ASR and SLU. SILENCE selectively mask the short-dependent sensitive words
309 while retaining the long-dependent SLU intents. Together with the differentiable mask generator,
310 SILENCE shows superior end-to-end inference speedup and privacy protection under different attack
311 scenarios.

312 **Limitations:** While for the first time, SILENCE provides a feasible privacy-preserving solution for
313 wimpy audio devices, it introduces a huge design space for mask generator structures. The mask
314 generator is akin to a lock; a genius lock design can protect privacy in the smallest of spaces, but a
315 poor lock design can be bulky and easily broken. In this work, we simply inherit the SLU model
316 structure and instantiate three sub-models from it to demonstrate better efficiency than previous
317 encoders. Researchers can explore other structures for a better privacy-accuracy-efficiency trade-
318 off. We will open-source all the code and checkpoints to facilitate further research in this direction.

319 Some other potential limitations about lossy privacy-preserving capacity, the need for fine-tuning
320 the cloud SLU model and the scope of defended threat model are thoroughly discussed below for
321 further clarification.

322 **Is current privacy-preserving capacity enough?** The quantitative WER 80% is considered secure
323 enough, as previous encoders have strived to reach that level [44, 10]. And some SLU transcripts
324 contain the intent word, so the successfully inferred word might be a non-private intent word. For
325 instance, in one test audio transcript, “I want some jazz music to play”, the intent is ‘scenario’:
326 ‘play’, ‘action’: ‘music’. The interpretation of the malicious cloud ASR, “all subjects were used
327 to play”, is acceptable since the predicted phrase “to play” contains no private information. This
328 scenario is typical for most audios; we managed to preserve 90% of the private entities in Figure 6.
329 This achievement matches the SoTA in privacy-preserving capacity, with up to $30\times$ lower latency
330 and $100\times$ memory reduction.

331 **Why and how to fine-tune the cloud SLU Model?** Initially, the cloud SLU is a generic pre-trained
332 speech model lacking the capability to accurately understand personalized user intent. It is crucial
333 to fine-tune the cloud SLU for better personalized intent understanding³. Secondly, while short-
334 dependent masking does not eliminate intent information, it does impact specific details within the
335 attention map, as depicted in Figure 4(b). Fine-tuning the cloud SLU model helps mitigate this
336 impact and enhances the understanding of the user’s intent.

337 Currently, cloud service providers have already offered APIs that allow users to fine-tune their per-
338 sonalized cloud speech model. For example, Azure has introduced the Custom Speech service [8],
339 which enables users to fine-tune the model for improved personalized outcomes. In this work, we
340 simulate the tunable cloud model using the open-source model to perform more detailed analysis,
341 such as different attacking scenarios

342 **Could private semantic detection attack be prevented?** SILENCE does not initially target private
343 semantic detection attacks. For example, eavesdropping on specific financial words and political
344 framing are *out-of-scope*. However, we can offer defense capabilities against them as discussed
345 below. The mask generator, controlled by the user, is trained to scrub utterances unrelated to the
346 public intent. Private entities not predefined by the user are almost never included in the masked
347 audio. Therefore, even if an attacker possesses a well-defined semantic and the mask generator,
348 training the detection threat model is challenging because the synthetic masked audio lacks clear
349 representations of the private semantic. Consequently, though not initially designed for this purpose,
350 our mask generators successfully discourage the malicious cloud provider from detecting private
351 semantics.

³Note that a general speech model is sufficient for training the local mask generator in Figure 5 step (1a), as the focus is not on generating precise intent but rather on obtaining a coarse-grained distribution of numerical logits to facilitate mask generator training.

References

- 352
- 353 [1] <https://openai.com/blog/chatgpt-can-now-see-hear-and-speak>.
- 354 [2] <https://huggingface.co/models?sort=downloads>.
- 355 [3] <https://safeatlast.co/blog/siri-statistics/>.
- 356 [4] [https://www.cnbc.com/2019/08/28/apple-apologizes-for-listening-to-siri-](https://www.cnbc.com/2019/08/28/apple-apologizes-for-listening-to-siri-conversations.html)
357 [conversations.html](https://www.cnbc.com/2019/08/28/apple-apologizes-for-listening-to-siri-conversations.html).
- 358 [5] [https://www.st.com/en/microcontrollers-microprocessors/stm32h7-](https://www.st.com/en/microcontrollers-microprocessors/stm32h7-series.html)
359 [series.html](https://www.st.com/en/microcontrollers-microprocessors/stm32h7-series.html).
- 360 [6] <https://www.raspberrypi.com/products/raspberry-pi-4-model-b/>.
- 361 [7] <https://developer.nvidia.com/embedded/jetson-tx2>.
- 362 [8] [https://azure.microsoft.com/en-us/blog/improve-speech-to-text-accuracy-](https://azure.microsoft.com/en-us/blog/improve-speech-to-text-accuracy-with-azure-custom-speech/)
363 [with-azure-custom-speech/](https://azure.microsoft.com/en-us/blog/improve-speech-to-text-accuracy-with-azure-custom-speech/).
- 364 [9] Shima Ahmed, Amrita Roy Chowdhury, Kassem Fawaz, and Parmesh Ramanathan. *Pr ϵ ech*:
365 A system for privacy-preserving speech transcription. [arXiv preprint arXiv:1909.04198](https://arxiv.org/abs/1909.04198) v2,
366 2019.
- 367 [10] Ranya Aloufi, Hamed Haddadi, and David Boyle. Privacy-preserving voice analysis via dis-
368 entangled representations. In *Proceedings of the 2020 ACM SIGSAC Conference on Cloud*
369 *Computing Security Workshop*, pages 1–14, 2020.
- 370 [11] Siddhant Arora, Siddharth Dalmia, Xuankai Chang, Brian Yan, Alan Black, and Shinji Watan-
371 abe. Two-pass low latency end-to-end spoken language understanding. [arXiv preprint](https://arxiv.org/abs/2207.06670)
372 [arXiv:2207.06670](https://arxiv.org/abs/2207.06670), 2022.
- 373 [12] Alexei Baeovski, Yuhao Zhou, Abdelrahman Mohamed, and Michael Auli. *wav2vec 2.0*:
374 A framework for self-supervised learning of speech representations. *Advances in neural*
375 *information processing systems*, 33:12449–12460, 2020.
- 376 [13] Emanuele Bastianelli, Andrea Vanzo, Pawel Swietojanski, and Verena Rieser. Slurp: A spoken
377 language understanding resource package. [arXiv preprint arXiv:2011.13205](https://arxiv.org/abs/2011.13205), 2020.
- 378 [14] Jasmijn Bastings, Wilker Aziz, and Ivan Titov. Interpretable neural predictions with differen-
379 tiable binary variables. [arXiv preprint arXiv:1905.08160](https://arxiv.org/abs/1905.08160), 2019.
- 380 [15] Yike Chen, Ming Gao, Yimin Li, Lingfeng Zhang, Li Lu, Feng Lin, Jinsong Han, and Kui Ren.
381 Big brother is listening: An evaluation framework on ultrasonic microphone jammers. In *IEEE*
382 *INFOCOM 2022-IEEE Conference on Computer Communications*, pages 1119–1128. IEEE,
383 2022.
- 384 [16] Peng Cheng and Utz Roedig. Personal voice assistant security and privacy—a survey.
385 *Proceedings of the IEEE*, 110(4):476–507, 2022.
- 386 [17] Leigh Clark, Philip Doyle, Diego Garaialde, Emer Gilmartin, Stephan Schlögl, Jens Edlund,
387 Matthew Aylett, João Cabral, Cosmin Munteanu, Justin Edwards, et al. The state of speech in
388 hci: Trends, themes and challenges. *Interacting with computers*, 31(4):349–371, 2019.
- 389 [18] Trung Dang, Om Thakkar, Swaroop Ramaswamy, Rajiv Mathews, Peter Chin, and Françoise
390 Beaufays. A method to reveal speaker identity in distributed asr training, and how to counter
391 it. In *ICASSP 2022-2022 IEEE International Conference on Acoustics, Speech and Signal*
392 *Processing (ICASSP)*, pages 4338–4342. IEEE, 2022.
- 393 [19] Nicola De Cao, Michael Schlichtkrull, Wilker Aziz, and Ivan Titov. How do decisions emerge
394 across layers in neural models? interpretation with differentiable masking. [arXiv preprint](https://arxiv.org/abs/2004.14992)
395 [arXiv:2004.14992](https://arxiv.org/abs/2004.14992), 2020.

- 396 [20] Keqi Deng, Songjun Cao, Yike Zhang, and Long Ma. Improving hybrid ctc/attention end-to-
397 end speech recognition with pretrained acoustic and language models. In 2021 IEEE Automatic
398 Speech Recognition and Understanding Workshop (ASRU), pages 76–82. IEEE, 2021.
- 399 [21] Ye Dong, Wen-jie Lu, Yancheng Zheng, Haoqi Wu, Derun Zhao, Jin Tan, Zhicong Huang,
400 Cheng Hong, Tao Wei, and Wenguang Cheng. Puma: Secure inference of llama-7b in five
401 minutes. arXiv preprint arXiv:2307.12533, 2023.
- 402 [22] Lloyd E Emokpae, Roland N Emokpae, Wassila Lalouani, and Mohamed Younis. Smart mul-
403 timodal telehealth-iot system for covid-19 patients. IEEE Pervasive Computing, 20(2):73–80,
404 2021.
- 405 [23] Ming Gao, Yike Chen, Yajie Liu, Jie Xiong, Jinsong Han, and Kui Ren. Cancelling Speech
406 Signals for Speech Privacy Protection against Microphone Eavesdropping. Association for
407 Computing Machinery, New York, NY, USA, 2023.
- 408 [24] Oded Goldreich. Secure multi-party computation. Manuscript. Preliminary version,
409 78(110):1–108, 1998.
- 410 [25] Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil
411 Ozair, Aaron Courville, and Yoshua Bengio. Generative adversarial nets. Advances in neural
412 information processing systems, 27, 2014.
- 413 [26] Anmol Gulati, James Qin, Chung-Cheng Chiu, Niki Parmar, Yu Zhang, Jiahui Yu, Wei Han,
414 Shibo Wang, Zhengdong Zhang, Yonghui Wu, et al. Conformer: Convolution-augmented
415 transformer for speech recognition. arXiv preprint arXiv:2005.08100, 2020.
- 416 [27] Parisa Haghani, Arun Narayanan, Michiel Bacchiani, Galen Chuang, Neeraj Gaur, Pedro
417 Moreno, Rohit Prabhavalkar, Zhongdi Qu, and Austin Waters. From audio to semantics:
418 Approaches to end-to-end spoken language understanding. In 2018 IEEE Spoken Language
419 Technology Workshop (SLT), pages 720–726. IEEE, 2018.
- 420 [28] Diederik P Kingma and Max Welling. Auto-encoding variational bayes. arXiv preprint
421 arXiv:1312.6114, 2013.
- 422 [29] Naveen Kumar and Seul Chan Lee. Human-machine interface in smart factory: A systematic
423 literature review. Technological Forecasting and Social Change, 174:121284, 2022.
- 424 [30] Christos Louizos, Max Welling, and Diederik P Kingma. Learning sparse neural networks
425 through l_0 regularization. arXiv preprint arXiv:1712.01312, 2017.
- 426 [31] Microsoft. Azure asr. [https://azure.microsoft.com/en-us/products/ai-services/
427 speech-to-text/](https://azure.microsoft.com/en-us/products/ai-services/speech-to-text/).
- 428 [32] Nombulelo CC Noruwana, Pius Adewale Owolawi, and Temitope Mapayi. Interactive iot-
429 based speech-controlled home automation system. In 2020 2nd International Multidisciplinary
430 Information Technology and Engineering Conference (IMITEC), pages 1–8. IEEE, 2020.
- 431 [33] Vassil Panayotov, Guoguo Chen, Daniel Povey, and Sanjeev Khudanpur. Librispeech: an
432 asr corpus based on public domain audio books. In 2015 IEEE international conference on
433 acoustics, speech and signal processing (ICASSP), pages 5206–5210. IEEE, 2015.
- 434 [34] Cal Peyser, Ronny Huang Andrew Rosenberg Tara N Sainath, Michael Picheny, and
435 Kyunghyun Cho. Towards disentangled speech representations. arXiv preprint
436 arXiv:2208.13191, 2022.
- 437 [35] Jianwei Qian, Haohua Du, Jiahui Hou, Linlin Chen, Taeho Jung, and Xiang-Yang Li. Hide-
438 behind: Enjoy voice input with voiceprint unclonability and anonymity. In Proceedings of the
439 16th ACM Conference on Embedded Networked Sensor Systems, pages 82–94, 2018.
- 440 [36] Alec Radford, Jong Wook Kim, Tao Xu, Greg Brockman, Christine McLeavey, and Ilya
441 Sutskever. Robust speech recognition via large-scale weak supervision. In International
442 Conference on Machine Learning, pages 28492–28518. PMLR, 2023.

- 443 [37] Joel M Raja, Carol Elsakr, Sherif Roman, Brandon Cave, Issa Pour-Ghaz, Amit Nanda, Miguel
444 Maturana, and Rami N Khouzam. Apple watch, wearables, and heart rhythm: where do we
445 stand? Annals of translational medicine, 7(17), 2019.
- 446 [38] Mirco Ravanelli, Titouan Parcollet, Peter Plantinga, Aku Rouhe, Samuele Cornell, Loren
447 Lugosch, Cem Subakan, Nauman Dawalatabad, Abdelwahab Heba, Jianyuan Zhong, et al.
448 Speechbrain: A general-purpose speech toolkit. arXiv preprint arXiv:2106.04624, 2021.
- 449 [39] Subendhu Rongali, Beiye Liu, Liwei Cai, Konstantine Arkoudas, Chengwei Su, and Wael
450 Hamza. Exploring transfer learning for end-to-end spoken language understanding. In
451 Proceedings of the AAAI Conference on Artificial Intelligence, volume 35, pages 13754–
452 13761, 2021.
- 453 [40] Steffen Schneider, Alexei Baevski, Ronan Collobert, and Michael Auli. wav2vec: Unsuper-
454 vised pre-training for speech recognition. arXiv preprint arXiv:1904.05862, 2019.
- 455 [41] Suranga Seneviratne, Yining Hu, Tham Nguyen, Guohao Lan, Sara Khalifa, Kanchana Thi-
456 lakarathna, Mahbub Hassan, and Aruna Seneviratne. A survey of wearable devices and chal-
457 lenges. IEEE Communications Surveys & Tutorials, 19(4):2573–2620, 2017.
- 458 [42] Ke Sun, Chen Chen, and Xinyu Zhang. " alexa, stop spying on me!" speech privacy protection
459 against voice assistants. In Proceedings of the 18th conference on embedded networked sensor
460 systems, pages 298–311, 2020.
- 461 [43] Rongxiang Wang and Felix Lin. Efficient deep speech understanding at the edge. arXiv preprint
462 arXiv:2311.17065, 2023.
- 463 [44] Yinggui Wang, Wei Huang, and Le Yang. Privacy-preserving end-to-end spoken language un-
464 derstanding. In Proceedings of the Thirty-Second International Joint Conference on Artificial
465 Intelligence, pages 5224–5232, 2023.
- 466 [45] Yingzhi Wang, Abdelmoumene Boumadane, and Abdelwahab Heba. A fine-tuned wav2vec
467 2.0/hubert benchmark for speech emotion recognition, speaker verification and spoken lan-
468 guage understanding. arXiv preprint arXiv:2111.02735, 2021.
- 469 [46] Thomas Wolf, Lysandre Debut, Victor Sanh, Julien Chaumond, Clement Delangue, Anthony
470 Moi, Pierric Cistac, Tim Rault, Rémi Louf, Morgan Funtowicz, et al. Huggingface’s trans-
471 formers: State-of-the-art natural language processing. arXiv preprint arXiv:1910.03771, 2019.
- 472 [47] Albert Zeyer, Ralf Schlüter, and Hermann Ney. Why does ctc result in peaky behavior? arXiv
473 preprint arXiv:2105.14849, 2021.
- 474 [48] Chengliang Zhang, Suyi Li, Junzhe Xia, Wei Wang, Feng Yan, and Yang Liu. Batchcrypt:
475 Efficient homomorphic encryption for cross-silo federated learning. In 2020 USENIX annual
476 technical conference (USENIX ATC 20), pages 493–506, 2020.

477 **NeurIPS Paper Checklist**

478 **(1) Claims**

479 Question: Do the main claims made in the abstract and introduction accurately reflect the
480 paper's contributions and scope?

481 Answer: [Yes]

482 Justification: Our contribution is outlined as a separated paragraph in §1.

483 Guidelines:

- 484 • The answer NA means that the abstract and introduction do not include the claims
485 made in the paper.
- 486 • The abstract and/or introduction should clearly state the claims made, including the
487 contributions made in the paper and important assumptions and limitations. A No or
488 NA answer to this question will not be perceived well by the reviewers.
- 489 • The claims made should match theoretical and experimental results, and reflect how
490 much the results can be expected to generalize to other settings.
- 491 • It is fine to include aspirational goals as motivation as long as it is clear that these
492 goals are not attained by the paper.

493 **(2) Limitations**

494 Question: Does the paper discuss the limitations of the work performed by the authors?

495 Answer: [Yes]

496 Justification: We thoroughly discuss the limitations of our work in §6.

497 Guidelines:

- 498 • The answer NA means that the paper has no limitation while the answer No means
499 that the paper has limitations, but those are not discussed in the paper.
- 500 • The authors are encouraged to create a separate "Limitations" section in their paper.
- 501 • The paper should point out any strong assumptions and how robust the results are to
502 violations of these assumptions (e.g., independence assumptions, noiseless settings,
503 model well-specification, asymptotic approximations only holding locally). The au-
504 thors should reflect on how these assumptions might be violated in practice and what
505 the implications would be.
- 506 • The authors should reflect on the scope of the claims made, e.g., if the approach was
507 only tested on a few datasets or with a few runs. In general, empirical results often
508 depend on implicit assumptions, which should be articulated.
- 509 • The authors should reflect on the factors that influence the performance of the ap-
510 proach. For example, a facial recognition algorithm may perform poorly when image
511 resolution is low or images are taken in low lighting. Or a speech-to-text system might
512 not be used reliably to provide closed captions for online lectures because it fails to
513 handle technical jargon.
- 514 • The authors should discuss the computational efficiency of the proposed algorithms
515 and how they scale with dataset size.
- 516 • If applicable, the authors should discuss possible limitations of their approach to ad-
517 dress problems of privacy and fairness.
- 518 • While the authors might fear that complete honesty about limitations might be used by
519 reviewers as grounds for rejection, a worse outcome might be that reviewers discover
520 limitations that aren't acknowledged in the paper. The authors should use their best
521 judgment and recognize that individual actions in favor of transparency play an impor-
522 tant role in developing norms that preserve the integrity of the community. Reviewers
523 will be specifically instructed to not penalize honesty concerning limitations.

524 **(3) Theory Assumptions and Proofs**

525 Question: For each theoretical result, does the paper provide the full set of assumptions and
526 a complete (and correct) proof?

527 Answer: [NA]

528 Justification: This paper does not include theoretical results.

529 Guidelines:

- 530 • The answer NA means that the paper does not include theoretical results.
- 531 • All the theorems, formulas, and proofs in the paper should be numbered and cross-
- 532 referenced.
- 533 • All assumptions should be clearly stated or referenced in the statement of any theo-
- 534 rems.
- 535 • The proofs can either appear in the main paper or the supplemental material, but if
- 536 they appear in the supplemental material, the authors are encouraged to provide a
- 537 short proof sketch to provide intuition.
- 538 • Inversely, any informal proof provided in the core of the paper should be comple-
- 539 mented by formal proofs provided in appendix or supplemental material.
- 540 • Theorems and Lemmas that the proof relies upon should be properly referenced.

541 (4) Experimental Result Reproducibility

542 Question: Does the paper fully disclose all the information needed to reproduce the main
543 experimental results of the paper to the extent that it affects the main claims and/or conclu-
544 sions of the paper (regardless of whether the code and data are provided or not)?

545 Answer: [\[Yes\]](#)

546 Justification: We provide detailed instructions on how to reproduce the main experimental
547 results in §4. We will open-source the code and data upon acceptance.

548 Guidelines:

- 549 • The answer NA means that the paper does not include experiments.
- 550 • If the paper includes experiments, a No answer to this question will not be perceived
- 551 well by the reviewers: Making the paper reproducible is important, regardless of
- 552 whether the code and data are provided or not.
- 553 • If the contribution is a dataset and/or model, the authors should describe the steps
- 554 taken to make their results reproducible or verifiable.
- 555 • Depending on the contribution, reproducibility can be accomplished in various ways.
- 556 For example, if the contribution is a novel architecture, describing the architecture
- 557 fully might suffice, or if the contribution is a specific model and empirical evaluation,
- 558 it may be necessary to either make it possible for others to replicate the model with
- 559 the same dataset, or provide access to the model. In general, releasing code and data
- 560 is often one good way to accomplish this, but reproducibility can also be provided via
- 561 detailed instructions for how to replicate the results, access to a hosted model (e.g., in
- 562 the case of a large language model), releasing of a model checkpoint, or other means
- 563 that are appropriate to the research performed.
- 564 • While NeurIPS does not require releasing code, the conference does require all sub-
- 565 missions to provide some reasonable avenue for reproducibility, which may depend
- 566 on the nature of the contribution. For example
- 567 (a) If the contribution is primarily a new algorithm, the paper should make it clear
- 568 how to reproduce that algorithm.
- 569 (b) If the contribution is primarily a new model architecture, the paper should describe
- 570 the architecture clearly and fully.
- 571 (c) If the contribution is a new model (e.g., a large language model), then there should
- 572 either be a way to access this model for reproducing the results or a way to re-
- 573 produce the model (e.g., with an open-source dataset or instructions for how to
- 574 construct the dataset).
- 575 (d) We recognize that reproducibility may be tricky in some cases, in which case au-
- 576 thors are welcome to describe the particular way they provide for reproducibility.
- 577 In the case of closed-source models, it may be that access to the model is limited in
- 578 some way (e.g., to registered users), but it should be possible for other researchers
- 579 to have some path to reproducing or verifying the results.

580 (5) Open access to data and code

581 Question: Does the paper provide open access to the data and code, with sufficient instruc-
582 tions to faithfully reproduce the main experimental results, as described in supplemental
583 material?

584 Answer: [No]

585 Justification: We will open-source the code and data upon acceptance.

586 Guidelines:

- 587 • The answer NA means that paper does not include experiments requiring code.
- 588 • Please see the NeurIPS code and data submission guidelines ([https://nips.cc/
589 public/guides/CodeSubmissionPolicy](https://nips.cc/public/guides/CodeSubmissionPolicy)) for more details.
- 590 • While we encourage the release of code and data, we understand that this might not
591 be possible, so “No” is an acceptable answer. Papers cannot be rejected simply for not
592 including code, unless this is central to the contribution (e.g., for a new open-source
593 benchmark).
- 594 • The instructions should contain the exact command and environment needed to run to
595 reproduce the results. See the NeurIPS code and data submission guidelines ([https:
596 //nips.cc/public/guides/CodeSubmissionPolicy](https://nips.cc/public/guides/CodeSubmissionPolicy)) for more details.
- 597 • The authors should provide instructions on data access and preparation, including how
598 to access the raw data, preprocessed data, intermediate data, and generated data, etc.
- 599 • The authors should provide scripts to reproduce all experimental results for the new
600 proposed method and baselines. If only a subset of experiments are reproducible, they
601 should state which ones are omitted from the script and why.
- 602 • At submission time, to preserve anonymity, the authors should release anonymized
603 versions (if applicable).
- 604 • Providing as much information as possible in supplemental material (appended to the
605 paper) is recommended, but including URLs to data and code is permitted.

606 (6) Experimental Setting/Details

607 Question: Does the paper specify all the training and test details (e.g., data splits, hyper-
608 parameters, how they were chosen, type of optimizer, etc.) necessary to understand the
609 results?

610 Answer: [Yes]

611 Justification: We provide detailed instructions on how to reproduce the main experimental
612 results in §4.

613 Guidelines:

- 614 • The answer NA means that the paper does not include experiments.
- 615 • The experimental setting should be presented in the core of the paper to a level of
616 detail that is necessary to appreciate the results and make sense of them.
- 617 • The full details can be provided either with the code, in appendix, or as supplemental
618 material.

619 (7) Experiment Statistical Significance

620 Question: Does the paper report error bars suitably and correctly defined or other appropri-
621 ate information about the statistical significance of the experiments?

622 Answer: [No]

623 Justification: Error bars are not reported because of the time limit. We will attempt to add
624 them in the camera-ready version.

625 Guidelines:

- 626 • The answer NA means that the paper does not include experiments.
- 627 • The authors should answer "Yes" if the results are accompanied by error bars, confi-
628 dence intervals, or statistical significance tests, at least for the experiments that support
629 the main claims of the paper.
- 630 • The factors of variability that the error bars are capturing should be clearly stated (for
631 example, train/test split, initialization, random drawing of some parameter, or overall
632 run with given experimental conditions).

- 633 • The method for calculating the error bars should be explained (closed form formula,
634 call to a library function, bootstrap, etc.)
- 635 • The assumptions made should be given (e.g., Normally distributed errors).
- 636 • It should be clear whether the error bar is the standard deviation or the standard error
637 of the mean.
- 638 • It is OK to report 1-sigma error bars, but one should state it. The authors should prefer-
639 ably report a 2-sigma error bar than state that they have a 96% CI, if the hypothesis of
640 Normality of errors is not verified.
- 641 • For asymmetric distributions, the authors should be careful not to show in tables or
642 figures symmetric error bars that would yield results that are out of range (e.g. negative
643 error rates).
- 644 • If error bars are reported in tables or plots, The authors should explain in the text how
645 they were calculated and reference the corresponding figures or tables in the text.

646 **(8) Experiments Compute Resources**

647 Question: For each experiment, does the paper provide sufficient information on the com-
648 puter resources (type of compute workers, memory, time of execution) needed to reproduce
649 the experiments?

650 Answer: [Yes]

651 Justification: We provide detailed hardware information in §4 and the intended runtime in
652 §3.2 and §5.

653 Guidelines:

- 654 • The answer NA means that the paper does not include experiments.
- 655 • The paper should indicate the type of compute workers CPU or GPU, internal cluster,
656 or cloud provider, including relevant memory and storage.
- 657 • The paper should provide the amount of compute required for each of the individual
658 experimental runs as well as estimate the total compute.
- 659 • The paper should disclose whether the full research project required more compute
660 than the experiments reported in the paper (e.g., preliminary or failed experiments
661 that didn't make it into the paper).

662 **(9) Code Of Ethics**

663 Question: Does the research conducted in the paper conform, in every respect, with the
664 NeurIPS Code of Ethics <https://neurips.cc/public/EthicsGuidelines>?

665 Answer: [Yes]

666 Justification: We have reviewed the NeurIPS Code of Ethics and believe that our research
667 conforms to it.

668 Guidelines:

- 669 • The answer NA means that the authors have not reviewed the NeurIPS Code of Ethics.
- 670 • If the authors answer No, they should explain the special circumstances that require a
671 deviation from the Code of Ethics.
- 672 • The authors should make sure to preserve anonymity (e.g., if there is a special consid-
673 eration due to laws or regulations in their jurisdiction).

674 **(10) Broader Impacts**

675 Question: Does the paper discuss both potential positive societal impacts and negative
676 societal impacts of the work performed?

677 Answer: [Yes]

678 Justification: We have discussed and provided real-world examples of both positive and
679 negative societal impacts in §1 and §2.

680 Guidelines:

- 681 • The answer NA means that there is no societal impact of the work performed.
- 682 • If the authors answer NA or No, they should explain why their work has no societal
683 impact or why the paper does not address societal impact.

- 684
- 685
- 686
- 687
- 688
- 689
- 690
- 691
- 692
- 693
- 694
- 695
- 696
- 697
- 698
- 699
- 700
- 701
- 702
- Examples of negative societal impacts include potential malicious or unintended uses (e.g., disinformation, generating fake profiles, surveillance), fairness considerations (e.g., deployment of technologies that could make decisions that unfairly impact specific groups), privacy considerations, and security considerations.
 - The conference expects that many papers will be foundational research and not tied to particular applications, let alone deployments. However, if there is a direct path to any negative applications, the authors should point it out. For example, it is legitimate to point out that an improvement in the quality of generative models could be used to generate deepfakes for disinformation. On the other hand, it is not needed to point out that a generic algorithm for optimizing neural networks could enable people to train models that generate Deepfakes faster.
 - The authors should consider possible harms that could arise when the technology is being used as intended and functioning correctly, harms that could arise when the technology is being used as intended but gives incorrect results, and harms following from (intentional or unintentional) misuse of the technology.
 - If there are negative societal impacts, the authors could also discuss possible mitigation strategies (e.g., gated release of models, providing defenses in addition to attacks, mechanisms for monitoring misuse, mechanisms to monitor how a system learns from feedback over time, improving the efficiency and accessibility of ML).

703 **(11) Safeguards**

704 Question: Does the paper describe safeguards that have been put in place for responsible
705 release of data or models that have a high risk for misuse (e.g., pretrained language models,
706 image generators, or scraped datasets)?

707 Answer: [NA]

708 Justification: This paper is intended for privacy protection and does not involve high-risk
709 data or models.

710 Guidelines:

- 711
- 712
- 713
- 714
- 715
- 716
- 717
- 718
- 719
- 720
- The answer NA means that the paper poses no such risks.
 - Released models that have a high risk for misuse or dual-use should be released with necessary safeguards to allow for controlled use of the model, for example by requiring that users adhere to usage guidelines or restrictions to access the model or implementing safety filters.
 - Datasets that have been scraped from the Internet could pose safety risks. The authors should describe how they avoided releasing unsafe images.
 - We recognize that providing effective safeguards is challenging, and many papers do not require this, but we encourage authors to take this into account and make a best faith effort.

721 **(12) Licenses for existing assets**

722 Question: Are the creators or original owners of assets (e.g., code, data, models), used in
723 the paper, properly credited and are the license and terms of use explicitly mentioned and
724 properly respected?

725 Answer: [Yes]

726 Justification: We have properly cited the original code, data and models in §4.

727 Guidelines:

- 728
- 729
- 730
- 731
- 732
- 733
- 734
- The answer NA means that the paper does not use existing assets.
 - The authors should cite the original paper that produced the code package or dataset.
 - The authors should state which version of the asset is used and, if possible, include a URL.
 - The name of the license (e.g., CC-BY 4.0) should be included for each asset.
 - For scraped data from a particular source (e.g., website), the copyright and terms of service of that source should be provided.

- 735
- If assets are released, the license, copyright information, and terms of use in the package should be provided. For popular datasets, paperswithcode.com/datasets has curated licenses for some datasets. Their licensing guide can help determine the license of a dataset.
- 736
- For existing datasets that are re-packaged, both the original license and the license of the derived asset (if it has changed) should be provided.
- 737
- If this information is not available online, the authors are encouraged to reach out to the asset's creators.
- 738
- 739
- 740
- 741
- 742

743 **(13) New Assets**

744 Question: Are new assets introduced in the paper well documented and is the documenta-
745 tion provided alongside the assets?

746 Answer: [No]

747 Justification: We will provide detailed documentation for the new assets upon acceptance.

748 Guidelines:

- The answer NA means that the paper does not release new assets.
 - Researchers should communicate the details of the dataset/code/model as part of their submissions via structured templates. This includes details about training, license, limitations, etc.
 - The paper should discuss whether and how consent was obtained from people whose asset is used.
 - At submission time, remember to anonymize your assets (if applicable). You can either create an anonymized URL or include an anonymized zip file.
- 749
- 750
- 751
- 752
- 753
- 754
- 755
- 756

757 **(14) Crowdsourcing and Research with Human Subjects**

758 Question: For crowdsourcing experiments and research with human subjects, does the pa-
759 per include the full text of instructions given to participants and screenshots, if applicable,
760 as well as details about compensation (if any)?

761 Answer: [NA]

762 Justification: This paper does not involve crowdsourcing nor research with human subjects.

763 Guidelines:

- The answer NA means that the paper does not involve crowdsourcing nor research with human subjects.
 - Including this information in the supplemental material is fine, but if the main contribution of the paper involves human subjects, then as much detail as possible should be included in the main paper.
 - According to the NeurIPS Code of Ethics, workers involved in data collection, curation, or other labor should be paid at least the minimum wage in the country of the data collector.
- 764
- 765
- 766
- 767
- 768
- 769
- 770
- 771

772 **(15) Institutional Review Board (IRB) Approvals or Equivalent for Research with Human
773 Subjects**

774 Question: Does the paper describe potential risks incurred by study participants, whether
775 such risks were disclosed to the subjects, and whether Institutional Review Board (IRB)
776 approvals (or an equivalent approval/review based on the requirements of your country or
777 institution) were obtained?

778 Answer: [NA]

779 Justification: This paper does not involve crowdsourcing nor research with human subjects.

780 Guidelines:

- The answer NA means that the paper does not involve crowdsourcing nor research with human subjects.
 - Depending on the country in which research is conducted, IRB approval (or equivalent) may be required for any human subjects research. If you obtained IRB approval, you should clearly state this in the paper.
- 781
- 782
- 783
- 784
- 785

786
787
788
789
790

- We recognize that the procedures for this may vary significantly between institutions and locations, and we expect authors to adhere to the NeurIPS Code of Ethics and the guidelines for their institution.
- For initial submissions, do not include any information that would break anonymity (if applicable), such as the institution conducting the review.

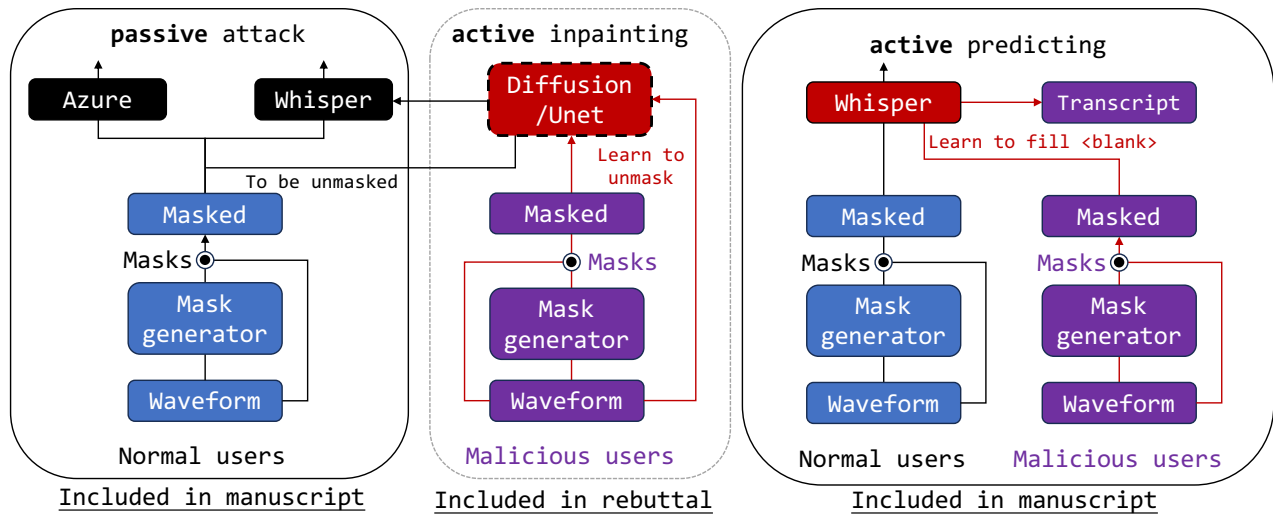


Figure 1: Mask generator and different attack scenarios, including both passive and active attacks.

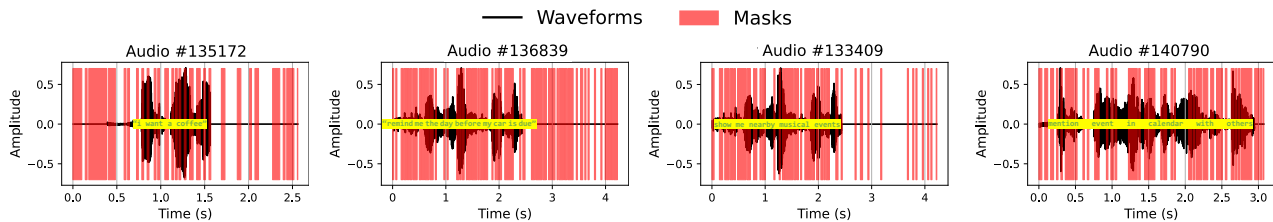


Figure 2: Illustration of the generated masks on audios selected randomly from SLURP. Local utterances are efficiently disrupted according to different transcripts patterns as highlighted within.

	PlainText	Azure	Naive Whisper	U-Net	CQT-Diff	Whisper predict (white box)
WER-SLU (%)	14.7	81.6	78.6	82.5	74.3	67.3
WER-ASR (%)	12.3	71.6	681.	71.4	65.9	64.4

Table 1: Potential attack Word Error Rate (WER) under different attack scenarios.

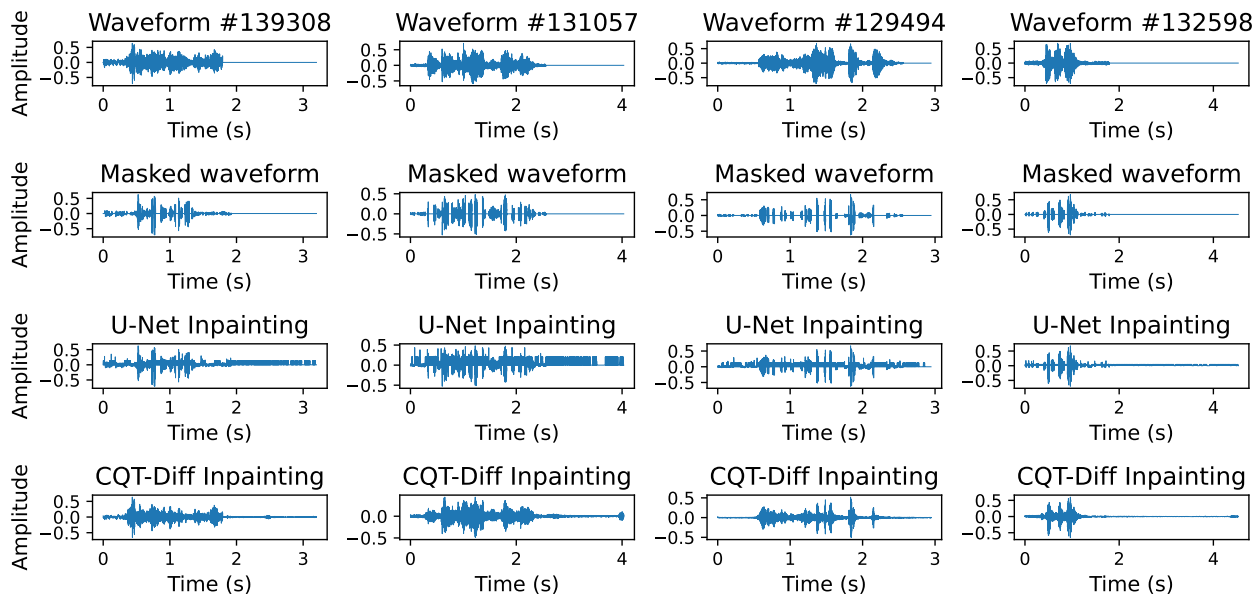


Figure 3: The reconstructed waveforms of different active inpainting attacks. Dataset: SLURP.