# PRIVACY-PRESERVING MECHANISMS ENABLE CHEAP VERIFIABLE INFERENCE OF LLMS

**Anonymous authors**Paper under double-blind review

000

001

002003004

010 011

012

013

014

015

016

017

018

019

021

025

026

027 028 029

031

032

033

034

037

040

041

042

043

044

046

047

048

051

052

#### **ABSTRACT**

As large language models (LLMs) continue to grow in size, fewer users are able to host and run models locally. This has led to increased use of third-party hosting services. However, in this setting, there is a lack of guarantees on the computation performed by the inference provider. For example, a dishonest provider may replace an expensive large model with a cheaper-to-run weaker model and return the results from the weaker model to the user. Existing tools to verify inference typically rely on methods from cryptography such as zero-knowledge proofs (ZKPs), but these add significant computational overhead, and remain infeasible for use for large models. In this work, we develop a new insight – that given a method for performing private LLM inference, one can obtain forms of verified inference at marginal extra cost. Specifically, we propose three new protocols, each of which leverage privacy-preserving LLM inference in order to provide different guarantees over the inference that was carried out. Our approaches are cheap, requiring the addition of a few extra tokens of computation, and have little to no downstream impact. As the fastest privacy-preserving inference methods are typically faster than ZK methods, the proposed protocols also improve verification runtime. Our work provides novel insights into the connections between privacy and verifiability in LLM inference.

# 1 Introduction

Large language models (LLMs) have increased significantly in size over the last few years. Recent models achieving cutting-edge performance (DeepSeek-AI et al., 2025; Qwen et al., 2025; Team et al., 2025), for example, now often contain hundreds of billions of parameters. The hardware requirements to run these models are often too high for individuals, or even organizations, to run on their own, leading to a significant growth in demand for third-party LLM inference providers. However, this trend raises critical concerns about the integrity and trustworthiness of the services provided, particularly in the growing decentralized inference space. In this setting, any entity with surplus computational resources can offer to complete computational tasks, such as LLM inference, for another user. As the providers in this setting are often not subject to strict vetting, it is imperative to ensure that the service paid for is actually one that is performed by the provider.

Traditionally, the verification of outsourced computation has been addressed through cryptographic methods, such as zero-knowledge proofs (ZKPs). Although offering strong theoretical guarantees, these methods often introduce substantial computational overhead for either the prover (the inference provider) or the verifier (the user), or both. Despite significant progress in recent years, the state-of-the-art for ZK verification of LLM inference remains thousands of times slower than vanilla inference (Sun et al., 2024), rendering it infeasible for large models, which are particularly likely to be in demand for third-party inference provision.

A related concern for third-party compute provision is that of *privacy-preservation*. Performing LLM inference for another party requires the user to share their prompts, resulting in a loss of privacy. Therefore, a seemingly orthogonal line of work in recent years has focused on privacy-preserving computation. These include methods such as secure multi-party computation (SMPC) and fully homomorphic encryption (FHE).

Our work examines the question: if a privacy mechanism is already in use, can this be leveraged to provide verification of the LLM inference computation as well? We answer this question

in the affirmative; specifically, we propose three simple but novel protocols, 'logit fingerprinting', 'logit fingerprinting with noise', and 'key appending', that use privacy to obtain differing levels of verification guarantees. We examine the costs and security properties of each of these protocols. Although our protocols have limitations and do not offer identical guarantees to those of ZK, we show that they are robust to many varieties of attacks. Moreover, we demonstrate that our 'logit fingerprinting with noise' protocol run with an SMPC method, SIGMA (Gupta et al., 2023), is  $\sim 15\times$  faster than the state-of-the-art ZK method for proof of LLM inference on a single forward pass of Llama-2-7B (Touvron et al., 2023). We further hope that connecting privacy and verification for LLM inference will spur the creation of improved protocols and further research in this area.

# 2 BACKGROUND & RELATED WORK

**Privacy-preserving inference.** There are four main families of privacy-preserving LLM inference methods that have been proposed in the literature: SMPC (Secure Multi-Party Computation), FHE (Fully Homomorphic Encryption), TEEs (Trusted Execution Environments), and statistical methods. SMPC and FHE are general privacy-preserving computation methods which provide strong guarantees on computational indistinguishability of the inputs. Both methods add significant overhead to plaintext computation; for SMPC, a large component of this is communication between the multiple parties involved. Recently, both SMPC (Huang et al., 2022; Hao et al., 2022; Pang et al., 2023; Akimoto et al., 2023; Dong et al., 2023; Li et al., 2023) and FHE (Moon et al., 2024; Zhang et al., 2024) have been applied to LLM inference. Our protocols are agnostic to the exact method used, though differing attacks are possible with each choice of privacy mechanism. For a more detailed account of these methods, see Section A.1.

**Verifiable inference.** Zero-knowledge proofs (ZKPs) are a class of methods that allows one party, the prover, to prove to another party, the verifier, that a staement is true, without revealing any additional information beyond the proof itself. ZK methods have recently been applied to proving LLM inference, such as in Sun et al. (2024); Qu et al. (2025). However, these approaches have significant overhead, and remain thousands of times slower than vanilla inference. By contrast, recent work has introduced 'statistical' methods for verifiable LLM inference, where the guarantees are relaxed in order to reduce the overhead added. For a more detailed account of these methods, see Section A.2.

Connections between privacy and verification. The connection between privacy and verification has not been extensively studied previously. Perhaps the closest work is MPC-in-the-Head (Ishai et al., 2007), which introduced a zero-knowledge verification protocol by utilizing any SMPC protocol. The protocol comes with steep costs for both the prover and verifier. For example, the prover must not only locally simulate every party in the underlying MPC execution but also repeat the computation multiple times. On the verifier's end, the party must perform several confirmation tasks, including recomputing opened views, consistency checks, and typically engage in multiple rounds of checking to achieve acceptable soundness. The crucial distinction of our suggested protocols to MPC-in-the-Head is that we use the privacy scheme *directly* to encode inexpensive secrets that are easily verifiable. To the best of our knowledge, there has not previously been any work that specifically examines the relationship between privacy-preserving LLM inference and verifiable inference of LLMs in this way.

#### 3 THREAT MODEL

We consider a setting with two primary roles: the **user**, who also acts as the verifier, and an **inference provider**, who also acts as the prover. The user wishes to run inference with a model M on their prompts x, but cannot do so themselves due to e.g. lack of computational resources. They therefore request the inference provider to perform inference on x with M. The inference provider is untrusted and may act as an adversary without behavorial constraints; other external adversaries are out of scope. We assume the use of a privacy-preserving mechanism providing computational indistinguishability of the inputs to ensure that the inference provider cannot view x. The model weights are assumed to be public. Our security goal is verifiability – that is, ensuring that the output the inference provider returns to the user can be verified as being the correct forward pass on the requested model on the given privatized prompt.

# 4 PROTOCOL 1: LOGIT FINGERPRINTING

Our first proposal for obtaining verification cheaply given access to a privacy-preserving method of LLM inference is **logit fingerprinting**. We hypothesize that the logit vector returned by performing a forward pass on any set of tokens on modern LLMs is a highly unique 'fingerprint' of the model. Our proposed protocol leverages this property to provide inference verification as follows:

- 1. First, the user inserts K sentinel tokens into the tokenized prompt, at random positions within the prompt. Call these positions  $p_1, p_2, ..., p_K$ . These K tokens are taken randomly from a public cache C, consisting of many such length K sequences.
- Next, the user creates the 2D attention mask to be used by the LLM by taking their desired attention mask (e.g., lower triangular for decoder-only LLMs) and inserting rows and columns as follows.
  - Add a row at  $p_i$  that is 0 everywhere except positions  $p_j \forall j \leq i$ , where it is set to 1.
  - Add a column at  $p_i$  that is 0 everywhere except positions  $p_i \, \forall \, j \geq i$ , where it is set to 1.
- 3. The attention mask and augmented tokenized prompt are given to the inference provider under a privacy-preserving scheme, and the inference provider carries out a forward pass, and returns the output logit vector at all token positions to the user.
- 4. The user verifies that the sentinel token logits match against the precomputed cached logits for that specific model.

The construction of the attention mask is such that the sentinel tokens do not attend to, and are not attended by, any of the original prompt tokens, but they do attend to each other in standard autoregressive fashion. This also ensures that sentinel tokens have no downstream impact on the original prompt when inference is performed. A formal description of this procedure is given in Section B.

#### 4.1 Cost Analysis

**Inference provider (prover).** Excluding the overhead of the private inference scheme, the total number of extra operations is a factor of  $\frac{K}{N}$ , where N is the length of the original prompt. As we discuss in Section 4.2, K can be set to be as small as 3 and retain strong security properties, so this is very small for reasonably sized N. Furthermore, if the privacy scheme supports parallelized inference, this can result in almost no extra runtime.

**User (verifier).** The verifier is required to pick a sequence from a public cache and perform a matching on the returned logits against the same cache. The cost of this is minimal and does not require specialized hardware.

Construction of the cache. Constructing the cache entails an initial computational cost and also must be performed by a trusted party, since it underpins the correctness of the protocol. Ideally, this responsibility is delegated to an entity with sufficient computational resources to produce a verifiable proof of correctness, for example, in the form of a zero-knowledge proof. Although the computational expense of this might be significant, the cost is incurred only once and is then amortized across all subsequent inference calls, including potentially all prover-verifier pairs.

#### 4.2 SECURITY ANALYSIS

In this section, we assume that logits are indeed unique fingerprints of models. We perform analysis across a range of models in Section 4.3 to verify this is the case.

In order for the inference provider to not be able to guess the logits to return for the sentinel tokens, the set of sentinel tokens must be randomly chosen from a large set of possibilities. The crux of this protocol is that the inference provider cannot determine which of the possibilities is specifically being asked for in any particular instance due to the privacy mechanism.

**Probabilistic attacks.** This protocol utilizes two elements of randomization: the choice of the sentinel tokens, and their positions. For the former, if the user selects the sequence uniformly at random from a cache of size |C|, then a dishonest inference provider can guess it with probability 1/|C|. |C| can

therefore be set to desired tolerances. For the latter, under a privacy-preserving mechanism that also preserves tensor structure (such as SMPC), correctly guessing the sentinel tokens' exact positions is sufficient for a successful attack: the inference provider can perform the forward pass on only those components, and return arbitrary values for the other token positions. However, this occurs with probability  $\binom{N+K}{K}^{-1}$ , where N is the length of the original prompt. When K=3, for example, with N=14, this is less than  $1\mathrm{e}{-3}$ , and it drops further with increasing N=100 to circa  $1\mathrm{e}{-6}$ .

A related attack is to perform computation only on a random subset of the token indices. In the most extreme case, a dishonest provider takes N+K-1 tokens, i.e. excludes exactly one token. The probability that all sentinel tokens are still selected (hence successfully passing verification) is  $\frac{N}{N+K}$ , requiring an infeasibly large K to make secure – although it should be noted that in this case the dishonest provider is saving very little computation over honest behavior.

**Approximation attacks.** Another line of possible attacks are attempts by the inference provider to use a different model – especially, cheaper-to-run replacements – that still succeed in passing verification. Such alternatives could include smaller models from the same model family or approximations to the models by using e.g. low-rank projections of the weights. We perform experiments to test the robustness of the protocol to each of the above in Section 4.3 and find that verification fails immediately when any of the above are attempted.

#### 4.3 EXPERIMENTS

**Setup.** We test the claim from Section 4.2 that pre–softmax logits can serve as model fingerprints. For each model m, we sample  $N=50{,}000$  token sequences of fixed length K=3 from the model's token vocabulary (excluding special tokens). Given a sequence  $\mathbf{t}=(t_1,t_2,t_3)$ , we run a forward pass and record the next-token *logit vectors* at each position,  $\ell_m^{(k)}(\mathbf{t}) \in \mathbb{R}^{V_m}$  for  $k \in \{1,2,3\}$ , where  $V_m$  is the vocabulary size of model m. We define the *logit fingerprint* 

$$\phi_m(\mathbf{t}) = \operatorname{concat}(\ell_m^{(1)}(\mathbf{t}), \ell_m^{(2)}(\mathbf{t}), \ell_m^{(3)}(\mathbf{t})) \in \mathbb{R}^{3V_m},$$

and compare fingerprints using L1 distance. We test on Llama 3.2 Instruct 1B, 3B, and 8B (Grattafiori et al., 2024), and on Qwen 2.5 Instruct 0.5B, 1.5B, 3B and 7B (Qwen et al., 2025). Comparisons are performed on FP32 logits; dropout is disabled.

#### 4.3.1 Honest Behaviors

**Floating point non-determinism.** We first provide context on the expected L1 distance due to non-determinism of floating-point operations (Shanmugavelu et al., 2024). This can be constituted as honest behavior; although, it is possible to also require an exact match, which would entail detecting hardware and batched-inference deviations. We run the *same* sequence multiple times with different batch sizes on GPU to measure this. We observe a maximum L1 deviation in doing so across all models tested of 10.90.

#### 4.3.2 DISHONEST BEHAVIORS

We now test a wide range of dishonest behaviors and strategies.

**Intra-model.** Within each model, we compute the nearest-neighbor similarity among fingerprints from *distinct* sequences (i.e.  $t \neq s$ ). Across N = 50k samples per model, there are no exact matches; the closest pair has an L1 distance of 2909.

**Within-family.** For the Llama family, the smallest L1 distance of logits we obtain is 329096. For the Qwen family, the minimum cross-model distance is 643719. These results indicate that even with a family of models, the logits are significantly different and suitable as fingerprints.

**Cross-family.** To enable comparisons across families with different vocabularies, we align dimensions by truncating the larger logit vectors to the smaller vocabulary size (i.e. comparing the first  $\min(V_m, V_{m'})$  coordinates). Under this conservative alignment, Llama–Qwen comparisons exhibit substantially higher distances than the within-family maxima reported above (qualitatively, well above 800000).

**Low-rank factorization.** We approximate the linear layers of Llama 3.2 1B Instruct by replacing each weight matrix  $W \in \mathbb{R}^{d_{\text{in}} \times d_{\text{out}}}$  with a rank-r factorization  $W \approx UV^{\top}$ , where  $U \in \mathbb{R}^{d_{\text{in}} \times r}$  and  $V \in$ 

216	
217	
218	
219	
220	
221	
222	
223	
224	

Model Approximation	Llama	Qwen
Same model, GPU non-determinism	4.08	10.91
Same model, different sequence	2909	68326
Cross-Model, same sequence	329096	643719
Low-rank factorization $(r = 2047)$	833	7223
Low-rank factorization ( $r = 2040$ )	8029	124195
Low-rank factorization $(r = 2000)$	31194	216765
8-bit Quantization	3499	8746
Single step of finetuning	471	1090

Table 1: L1 distances of logit fingerprints across different experimental settings for the Llama and Qwen families of models. Honest behavior (GPU non-determinism) is clearly separated from dishonest behavior (all other rows).

 $\mathbb{R}^{d_{\text{out}} \times r}$ . The default hidden dimension of this model is 2048, so we test with  $r \in \{2047, 2040, 2000\}$ . Comparing fingerprints of 50k sequences between the full-rank and the low-rank variants, the minimum L1 distances observed are:

```
r = 2047:833.97, \quad r = 2040:8029.45, \quad r = 2000:31194.02.
```

Similarly, we obtained the following results for Qwen 2.5 3B Instruct:

```
r = 2047:7223, \quad r = 2040:124195, \quad r = 2000:216765.
```

**Quantization.** We next load Llama 3.2 1B Instruct and Qwen 2.5 3B Instruct in 8-bit precision using bitsandbytes and compare fingerprints to the full-precision (bfloat16) baseline. The minimum L1 distance is 3499 for Llama and 8746 for Qwen, again easily separated from the original model.

**Fine-tuning.** Finally, we evaluate robustness against model fine-tuning by comparing each of Llama 3.2 1B Instruct and Qwen 2.5 3B Instruct with a finetuned variant of each corresponding model on a single sample from FineWeb dataset for a single step. The minimum observed distance is 471 for Llama and 1090 for Qwen, consistent with the previous cases and again easily separable from the original model.

Our results are summarized in Table 1. The minimum L1 distance observed under dishonest behavior is 471, as seen in Llama's finetuning setting, while the maximum deviation with honest behavior due to floating point non-determinism is only 10.90 for Qwen. The significant difference allows for clear identification of honest vs. dishonest behavior; based on these results, we recommend using a matching threshold in the range of 15-20. Sequences whose logits differ by less than this threshold can be confidently regarded as originating from the same model; and even a single step of fine-tuning is easily detectable with this threshold.

# 4.4 LIMITATIONS

The main limitation of this protocol is that it can only be used to verify a single forward pass at a time, i.e. only generate a single new token, before requiring the user to repeat the protocol above with a fresh set of sentinel tokens and positions; otherwise, a dishonest provider could honestly perform the first forward pass (to pass verification) and provide spurious outputs for all subsequent forward passes. Thus, this protocol inherently requires user interaction for every step of token decoding. Another limitation is the vulnerability to the subsetting attack mentioned in Section 4.2. As such, we recommend that this protocol not be used in isolation with privacy mechanisms that retain tensor structure, such as SMPC methods.

# 5 PROTOCOL 2: LOGIT FINGERPRINTING WITH NOISE

The vulnerability of Protocol 1 to a subsetting attack reduces the space of privacy gadgets that it can be used with. Our second proposed protocol is designed to resist this attack. Our modification consists of adding randomly sampled noise to the token embeddings before they are passed into the LLM for the forward pass, and then using a lightweight predictor on the returned final hidden states to predict the noise that was used. Our proposed protocol is as follows:

- 270
- 271
- 272
- 274
- 275 276
- 277 278 279
- 281
- 284
- 287
- 288 289 290
- 291
- 293
- 295 296
- 297 298
- 299 300
- 301 302
- 303 304
- 305 306 307
- 308
- 310 311
- 312 313 314
- 315 316 317
- 318 319
- 320 321
- 322 323

- 1. First the user samples noise  $b \in \mathbb{R}_e^d$ , where  $d_e$  is the embedding dimension of the model being used for inference, from a discrete set of possibilities B.
- 2. The user concatenates the noise to the embedding of the original prompt  $e \in \mathbb{R}^{N \times d_e}$  in the  $d_e$  dimension, to obtain a tensor  $b_e \in \mathbb{R}^{N \times 2d_e}$ .
- 3. The user applies the previously trained *NoiseEmbedder* module on  $b_e$  to obtain  $e' \in \mathbb{R}^{N \times d_e}$ .
- 4. The user sends privatized e', augmented with K randomly-positioned sentinel tokens as in Protocol 1, to the inference provider.
- 5. The inference provider performs the forward pass and returns the final hidden states  $h \in$  $\mathbb{R}^{(N+K)\times d_h}$ , where  $d_h$  is the hidden dimension.
- 6. The user applies the logit-projection to the hidden states at the sentinel token positions, and checks the validity of these logits against the cache, as in Protocol 1.
- 7. The user applies the previously trained prediction module, the *NoisePredictor*, on h at the non-sentinel positions to obtain estimated b at each such position. If each obtained b matches the sampled b at that position, and the sentinel token logit check passes above, then the user can consider the inference to be verified.
- In the above procedure, the sentinel tokens are not modified by the sampled noise, and so can be compared against the cache as in Protocol 1. For the remainder of the tokens, the predicted noise is compared to the sampled noise to verify that the forward pass was indeed carried out on each token position. The complete procedure is formally described in Section C.

# 5.1 Cost Analysis

- **Inference provider (prover).** The cost to the inference provider is the same in this protocol as in
- User (verifier). In addition to the cost associated with the sentinel tokens, the user must now generate the sampled noise – which requires little computational cost – as well as run the NoiseEmbedder and NoisePredictor.
- **Construction of the cache.** This cost remains the same as in Protocol 1.
- Training of NoiseEmbedder and NoisePredictor. The NoiseEmbedder and NoisePredictor modules need to be trained for each different LLM in use. This entails an initial computational cost and also must be performed by a trusted party; however, similarly to the cache construction, this is a one-time cost that is then amortized over all subsequent inference calls on that model.
- In Section 5.3, we show that the NoiseEmbedder and NoisePredictor can be simple linear projections, so that the additional cost to the user and the training cost can be made low in practice.

# 5.2 SECURITY ANALYSIS

- We inherit the security analysis of Protocol 1 as it pertains to sentinel tokens that is, sentinel tokens remain effective markers of the model that was used for the forward pass and are able to detect even very close replacements. For the non-sentinel tokens, the crux of the protocol's security now rests on the predictability of the injected noise.
- Let the sample space size be given by |B|, and denote the accuracy of the prediction at token position n under honest inference by  $acc_n := P(b_n = b_n \mid honest)$ .
- **Honest provider (completeness).** If the inference provider is honest, the probability that the user incorrectly rejects the returned computation is given by the probability that there is at least one mismatch in the predicted noise:  $P(\text{incorrect rejection}) = 1 - \prod_{n=1}^{N} \text{acc}_n$ .
- **Dishonest provider (soundness).** If the inference provider is dishonest, the probability that the user incorrectly accepts the returned computation is given by the probability that  $b_n = b_n$  at all token positions n. Due to the privacy-preserving mechanism, the provider cannot know which  $b_n$  was used, so the probability that  $b_n = b_n$  for any particular n is upper bounded by  $\frac{1}{|B|}$ . In particular, the above

implies that in a leave-one-out subsetting attack as described in Section 4.2, the probability of success is at most  $\frac{1}{|B|}$ .

#### 5.3 EXPERIMENTS

In this section, we describe a performant and lightweight architecture of the NoiseEmbedder and NoisePredictor; and we demonstrate their performance on Llama-3.2-1B, evaluating on the FineWeb-Edu dataset (Lozhkov et al., 2024).

**NoiseEmbedder architecture.** This module consists of a learned embedding  $E \in \mathbb{R}^{|B| \times d_e}$ , and a linear layer that is applied to the concatenation of the learned noise embeddings and the original embedding, and produces a single combined embedding as an output. Therefore the linear layer has a weight matrix:  $W \in \mathbb{R}^{2d_e \times d_e}$ .

**NoisePredictor architecture.** This module consists of a linear layer that takes the final hidden layer representations from the forward pass of the LLM and outputs unnormalized logits over the sample space B. Therefore the linear layer has a weight matrix:  $W \in \mathbb{R}^{d_h \times |B|}$ .

We fine-tune the NoiseEmbedder and NoisePredictor modules **whilst keeping the original model weights frozen**. As we are adding noise to the model embeddings, we train to optimize for both the log-likelihood on the dataset, as well as the classification accuracy of the NoisePredictor, using the cross-entropy loss. Our training objective is therefore given by:

$$\mathcal{L}_{\theta,\phi} = \mathbb{E}_{x,y\sim\mathcal{D}} \mathbb{E}_{b\sim B} \Big[ -\log f \big( y \mid \text{NoiseEmbedder}_{\theta}(x,b) \big) + \\ \lambda \operatorname{CE} \big( \text{NoisePredictor}_{\phi} \big( f(\text{NoiseEmbedder}_{\theta}(x,b)) \big), b \big) \Big]$$

$$(1)$$

where x, y are the training data, f is the base model,  $\theta$  and  $\phi$  are the parameters of the NoiseEmbedder and NoisePredictor respectively, and  $\lambda$  is a hyperparameter to be tuned. In practice, we find best results applying the same sampled noise to every token in the sequence. This does not impact the security analysis of Section 5.2. For further training and hyperparameter details, see Section D.

**Results.** Despite using a very lightweight NoiseEmbedder and NoisePredictor, and not modifying the original model weights at all, we find that we are able to achieve  $\sim 99\%$  classification accuracy with |B|=100 without any worsening of the log-loss on the given dataset. In particular, the base model's log-prob is  $\sim 3.45$ , and we achieve a held-out evaluation set log-prob of  $\sim 3.43$  after training the modules, with noise injected.

#### 5.4 LIMITATIONS

In comparison with Protocol 1, this protocol is resistant to subset attacks, due to the introduction of noise at each token position. However, this protocol adds extra computational burden to the user – they must now perform additional NoiseEmbedder and NoisePredictor forward passes. Although we have shown that these can be effective even if comprising just a single linear layer each, there may be some cases where even this extra computational requirement cannot be met. Moreover, the user must now also perform projection of the final hidden states to the logits themselves, necessitating another matrix multiplication. There is also now the additional computational requirement of training the modules prior to deployment, in a trust-secured manner. Finally, although we are able to achieve good accuracy rates of 99% with |B|=100, we have neither perfect soundness nor completeness; we hope that future work is capable of improving on the results we present here.

#### 6 Protocol 3: Key Appending

A conceptually distinct proposal for obtaining verification cheaply under privacy assumptions is **key appending**. The high-level idea of this protocol is to ask the LLM to emit a randomly generated key at the end of its normal response – for example, 'strawberry reticent gestalt' – wrapped in a specific tag structure, and to verify at the end of inference that the key was correctly replicated.

**User prompt augmentation.** Given a user's original prompt p and a randomly sampled key  $w_1, w_2, \ldots, w_K$  of K words, p is augmented by appending the following instruction to the end:

```
At the end of your response, repeat this key: <key> w_1\ w_2\ \cdots\ w_K </key>\nDO NOT print anything else after it.
```

**Verification.** The encrypted response is returned to the user, who decrypts it and parses the contents of the  $\langle key \rangle$ ... $\langle key \rangle$  span. Verification succeeds if and only if the extracted string matches the originally sampled key. Using HTML-/XML-like tags facilitates robust parsing and prevents ambiguity in locating the verification key within the model's output.

Further details on system prompt modification and the stopping criterion are given in Section E.

The main benefit of this protocol as compared to Protocols 1 and 2 is that it does not require continual user interaction at every decoding step; it is an entirely non-interactive protocol.

#### 6.1 Cost Analysis

**Inference provider (Prover).** Adding an extra t tokens to the prompt adds an overhead of a factor of  $\frac{t}{N}$  operations (in addition to the system prompt and remaining augmentations, which are of constant length). Given that, for most tokenizers and English words, words are approximately 1–2 tokens in length, and that the protocol offers good security with just K=3 words (see Section E.2), this therefore introduces little extra overhead.

**User (Verifier).** Similar to the logit fingerprinting protocol, the verifier is required to perform minimal work. They must select a sequence of K words  $w_1w_2\cdots w_K$  and append them together with the augmentation template to the prompt, as well as prepend the system prompt. When the inference is complete, the verifier checks the words between the key tags of the decoded output against the original words  $w_1w_2\cdots w_K$ . Again, no specialized hardware is necessary.

#### 6.2 SECURITY ANALYSIS

**Probabilistic Attacks.** Suppose that tokenizing the K words results in a total of t tokens. An adversarial party must correctly guess each token exactly out of the total vocabulary, resulting in a success probability of  $\frac{1}{|V|^t}$ , where |V| is the vocabulary size. As modern LLMs typically have  $|V| \ge 1$ e5, with a key length of only 3 tokens, this is already on the order of 1e-15 or lower.

**Approximation Attacks.** This protocol is potentially vulnerable to the model approximation attacks as described in Section 4.2, especially if they largely retain the instruction following capabilities of the original model. We perform an in-depth examination of the viability of such in the specific case of use of SMPC for privacy preservation, with one honest party, in Section G. We find that in this setting, any such approximation does fail.

#### 6.3 EXPERIMENTS

We conduct two tests in this section. First, we test the ability of LLMs to perform the protocol successfully under honest behavior; this is analogous to the cryptographic property of *completeness*. Second, we test the performance impact of running the appending protocol on other downstream tasks. We find that models of large enough size  $\geq 3B$  parameters in scale, have nearly perfect key transcription rates. Further, we find minimal performance impact for models of size 8B or larger on downstream tasks. Further details of our results are given in Section E.

#### 6.4 LIMITATIONS

The main limitation of this protocol is that it cannot specify exactly the model that is being used; it only guarantees that the model is capable of performing the verification task. In the SMPC setting, if there is at least one honest participant, then we show in Section G that any approximations result in the verification failing. However, in the FHE or TEE setting, this remains a limitation. Furthermore, this protocol does not offer 100% completeness, although we see figures close to this (see Table 3).

## 7 OUTPERFORMING STATE-OF-THE-ART ZK INFERENCE

We have shown that privacy-preserving mechanisms can enable verified inference. In this section, we describe the performance of privacy-preserving inference and our protocols, compared to the standard approach of zero-knowledge (ZK) proofs of inference. In particular, we focus on secure multi-party computation (SMPC) and fully homomorphic encryption (FHE) schemes for the latter. For a detailed background on these methods, see Section A.1 and Section A.2.

The protocols we proposed in Section 4 and Section 5 are both compatible with FHE privacy schemes. However, state-of-the-art FHE schemes typically have greater overhead than ZK; for example, THOR (Moon et al., 2024) reports approximately 10 minutes for a single forward pass on an input of 128 tokens with BERT-Base (a model with 110M params), with GPU acceleration. By contrast, zkLLM reports just 74s of prover overhead for a forward pass with an input of 2048 tokens on OPT-125M. However, Protocol 2 (Section 5) is designed to resist tensor subset attacks, and is therefore also compatible for use with SMPC schemes. State-of-the-art SMPC schemes operate much faster than FHE.

We perform a direct comparison of our protocol with SMPC to zkLLM. We run zkLLM on Llama-2-7B (Touvron et al., 2023) on sequences of length 125, and measure the total prover time for a single complete forward pass and associated proof generation on a machine with an A6000 GPU. We compare this to the results described in SIGMA (Gupta et al., 2023). SIGMA is a 2-party SMPC scheme that is optimized for GPU acceleration. The authors of that paper report performance on a machine also accelerated with an A6000 GPU. We take the result from Table 5 of that paper indicating a total runtime of Llama-2-7B on SIGMA of 23s for a single forward pass, on sequence lengths of 128 tokens. The extra 3 tokens in the SIGMA setting correspond to the sentinel tokens used in our protocol.

Our results are shown in Table 2. We see that Protocol 2 under SIGMA is approximately  $\sim 15 \times$  faster than zkLLM.

Method	Time (s)
zkLLM	352
Protocol 2 w/ SIGMA (ours)	23

Table 2: Inference provider runtime for a single forward pass of Llama-2-7B on a prompt with length 125 tokens with state-of-the-art ZK method zkLLM, and with our Protocol 2 with SMPC protocol SIGMA as the privacy-preserving mechanism. SIGMA numbers are taken from Gupta et al. (2023).

**Discussion.** Our protocol as tested in a like-for-like setting is nearly  $15 \times$  faster than the state-of-the-art ZK method for proof of LLM inference. However, there are two key differences. First, ZK has fewer security assumptions. Although SMPC guarantees strong computational indistinguishability of its inputs in the non-colluding setting, it is vulnerable when all parties involved are dishonest and collude to pool their secret shares. By contrast, ZK is provably secure regardless of prover behavior assumptions. Second, our protocol still relies on statistical results, such as the accuracy of the NoisePredictor module. Therefore, our inference guarantees are not directly comparable to those produced by ZK methods.

Nevertheless, in settings where non-collusion can be ensured or encouraged, and where statistical guarantees are sufficient, our protocol offers a significant speedup over the state-of-the-art for proof of LLM inference.

#### 8 CONCLUSION

We have introduced three protocols for verifying LLM inference, given the use of privacy-preserving mechanisms. These protocols are cheap for both the prover and the verifier and have little to no downstream impact. Future work may focus on mitigating the limitations of our protocols, for example by (1) boosting efficiency during many-token generation, (2) improving the statistical guarantees, or (3) guaranteeing resistance to attacks. We believe that connecting privacy and verifiability, particularly in LLM inference, will inspire future work on new and improved protocols.

# REFERENCES

486

487

488

489

490 491

492

493

494

495

496

497

498

499

500

501

504

505

506

507

508

509

510

511

512

513

514

515

516

517

519

520

521

522

523

524 525

526

527 528

529

530

531

532

533

534

535

536

538

- Yoshimasa Akimoto, Kazuto Fukuchi, Youhei Akimoto, and Jun Sakuma. Privformer: Privacy-preserving transformer with mpc. In 2023 IEEE 8th European Symposium on Security and Privacy (EuroS&P), pp. 392–410, 2023. doi: 10.1109/EuroSP57164.2023.00031.
- Argilla. argilla/databricks-dolly-15k-curated-en [dataset]. https://huggingface.co/datasets/argilla/databricks-dolly-15k-curated-en, 2023. Downloaded from Hugging Face Hub on 2025-08-29.
- Jung Hee Cheon, Andrey Kim, Miran Kim, and Yongsoo Song. Homomorphic encryption for arithmetic of approximate numbers. In *International conference on the theory and application of* cryptology and information security, pp. 409–437. Springer, 2017.
- DeepSeek-AI, Daya Guo, Dejian Yang, Haowei Zhang, Junxiao Song, Ruoyu Zhang, Runxin Xu, Qihao Zhu, Shirong Ma, Peiyi Wang, Xiao Bi, Xiaokang Zhang, Xingkai Yu, Yu Wu, Z. F. Wu, Zhibin Gou, Zhihong Shao, Zhuoshu Li, Ziyi Gao, Aixin Liu, Bing Xue, Bingxuan Wang, Bochao Wu, Bei Feng, Chengda Lu, Chenggang Zhao, Chengqi Deng, Chenyu Zhang, Chong Ruan, Damai Dai, Deli Chen, Dongjie Ji, Erhang Li, Fangyun Lin, Fucong Dai, Fuli Luo, Guangbo Hao, Guanting Chen, Guowei Li, H. Zhang, Han Bao, Hanwei Xu, Haocheng Wang, Honghui Ding, Huajian Xin, Huazuo Gao, Hui Qu, Hui Li, Jianzhong Guo, Jiashi Li, Jiawei Wang, Jingchang Chen, Jingyang Yuan, Junjie Qiu, Junlong Li, J. L. Cai, Jiaqi Ni, Jian Liang, Jin Chen, Kai Dong, Kai Hu, Kaige Gao, Kang Guan, Kexin Huang, Kuai Yu, Lean Wang, Lecong Zhang, Liang Zhao, Litong Wang, Liyue Zhang, Lei Xu, Leyi Xia, Mingchuan Zhang, Minghua Zhang, Minghui Tang, Meng Li, Miaojun Wang, Mingming Li, Ning Tian, Panpan Huang, Peng Zhang, Qiancheng Wang, Qinyu Chen, Qiushi Du, Ruiqi Ge, Ruisong Zhang, Ruizhe Pan, Runji Wang, R. J. Chen, R. L. Jin, Ruyi Chen, Shanghao Lu, Shangyan Zhou, Shanhuang Chen, Shengfeng Ye, Shiyu Wang, Shuiping Yu, Shunfeng Zhou, Shuting Pan, S. S. Li, Shuang Zhou, Shaoqing Wu, Shengfeng Ye, Tao Yun, Tian Pei, Tianyu Sun, T. Wang, Wangding Zeng, Wanjia Zhao, Wen Liu, Wenfeng Liang, Wenjun Gao, Wenqin Yu, Wentao Zhang, W. L. Xiao, Wei An, Xiaodong Liu, Xiaohan Wang, Xiaokang Chen, Xiaotao Nie, Xin Cheng, Xin Liu, Xin Xie, Xingchao Liu, Xinyu Yang, Xinyuan Li, Xuecheng Su, Xuheng Lin, X. Q. Li, Xiangyue Jin, Xiaojin Shen, Xiaosha Chen, Xiaowen Sun, Xiaoxiang Wang, Xinnan Song, Xinyi Zhou, Xianzu Wang, Xinxia Shan, Y. K. Li, Y. Q. Wang, Y. X. Wei, Yang Zhang, Yanhong Xu, Yao Li, Yao Zhao, Yaofeng Sun, Yaohui Wang, Yi Yu, Yichao Zhang, Yifan Shi, Yiliang Xiong, Ying He, Yishi Piao, Yisong Wang, Yixuan Tan, Yiyang Ma, Yiyuan Liu, Yongqiang Guo, Yuan Ou, Yuduan Wang, Yue Gong, Yuheng Zou, Yujia He, Yunfan Xiong, Yuxiang Luo, Yuxiang You, Yuxuan Liu, Yuyang Zhou, Y. X. Zhu, Yanhong Xu, Yanping Huang, Yaohui Li, Yi Zheng, Yuchen Zhu, Yunxian Ma, Ying Tang, Yukun Zha, Yuting Yan, Z. Z. Ren, Zehui Ren, Zhangli Sha, Zhe Fu, Zhean Xu, Zhenda Xie, Zhengyan Zhang, Zhewen Hao, Zhicheng Ma, Zhigang Yan, Zhiyu Wu, Zihui Gu, Zijia Zhu, Zijun Liu, Zilin Li, Ziwei Xie, Ziyang Song, Zizheng Pan, Zhen Huang, Zhipeng Xu, Zhongyu Zhang, and Zhen Zhang. Deepseek-r1: Incentivizing reasoning capability in llms via reinforcement learning, 2025. URL https://arxiv.org/abs/2501.12948.
- Ye Dong, Wen jie Lu, Yancheng Zheng, Haoqi Wu, Derun Zhao, Jin Tan, Zhicong Huang, Cheng Hong, Tao Wei, and Wenguang Chen. Puma: Secure inference of llama-7b in five minutes, 2023. URL https://arxiv.org/abs/2307.12533.
- Craig Gentry. Fully homomorphic encryption using ideal lattices. In *Proceedings of the Forty-First Annual ACM Symposium on Theory of Computing*, STOC '09, pp. 169–178, New York, NY, USA, 2009. Association for Computing Machinery. ISBN 9781605585062. doi: 10.1145/1536414. 1536440. URL https://doi.org/10.1145/1536414.1536440.
- O. Goldreich, S. Micali, and A. Wigderson. How to play any mental game. In *Proceedings of the Nineteenth Annual ACM Symposium on Theory of Computing*, STOC '87, pp. 218–229, New York, NY, USA, 1987. Association for Computing Machinery. ISBN 0897912217. doi: 10.1145/28395.28420. URL https://doi.org/10.1145/28395.28420.
- S Goldwasser, S Micali, and C Rackoff. The knowledge complexity of interactive proof-systems. In *Proceedings of the Seventeenth Annual ACM Symposium on Theory of Computing*, STOC '85, pp. 291–304, New York, NY, USA, 1985. Association for Computing Machinery. ISBN 0897911512. doi: 10.1145/22145.22178. URL https://doi.org/10.1145/22145.22178.

541

542

543

544

546

547

548

549

550

551

552

553

554

556

558

559

561

562

563

564

565

566

567

568

569

570

571

572

573

574

575

576

577

578

579

580

581

582

583

584

585

586

588

590

592

Aaron Grattafiori, Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Alex Vaughan, Amy Yang, Angela Fan, Anirudh Goyal, Anthony Hartshorn, Aobo Yang, Archi Mitra, Archie Sravankumar, Artem Korenev, Arthur Hinsvark, Arun Rao, Aston Zhang, Aurelien Rodriguez, Austen Gregerson, Ava Spataru, Baptiste Roziere, Bethany Biron, Binh Tang, Bobbie Chern, Charlotte Caucheteux, Chaya Nayak, Chloe Bi, Chris Marra, Chris McConnell, Christian Keller, Christophe Touret, Chunyang Wu, Corinne Wong, Cristian Canton Ferrer, Cyrus Nikolaidis, Damien Allonsius, Daniel Song, Danielle Pintz, Danny Livshits, Danny Wyatt, David Esiobu, Dhruv Choudhary, Dhruv Mahajan, Diego Garcia-Olano, Diego Perino, Dieuwke Hupkes, Egor Lakomkin, Ehab AlBadawy, Elina Lobanova, Emily Dinan, Eric Michael Smith, Filip Radenovic, Francisco Guzmán, Frank Zhang, Gabriel Synnaeve, Gabrielle Lee, Georgia Lewis Anderson, Govind Thattai, Graeme Nail, Gregoire Mialon, Guan Pang, Guillem Cucurell, Hailey Nguyen, Hannah Korevaar, Hu Xu, Hugo Touvron, Iliyan Zarov, Imanol Arrieta Ibarra, Isabel Kloumann, Ishan Misra, Ivan Evtimov, Jack Zhang, Jade Copet, Jaewon Lee, Jan Geffert, Jana Vranes, Jason Park, Jay Mahadeokar, Jeet Shah, Jelmer van der Linde, Jennifer Billock, Jenny Hong, Jenya Lee, Jeremy Fu, Jianfeng Chi, Jianyu Huang, Jiawen Liu, Jie Wang, Jiecao Yu, Joanna Bitton, Joe Spisak, Jongsoo Park, Joseph Rocca, Joshua Johnstun, Joshua Saxe, Junteng Jia, Kalyan Vasuden Alwala, Karthik Prasad, Kartikeya Upasani, Kate Plawiak, Ke Li, Kenneth Heafield, Kevin Stone, Khalid El-Arini, Krithika Iyer, Kshitiz Malik, Kuenley Chiu, Kunal Bhalla, Kushal Lakhotia, Lauren Rantala-Yeary, Laurens van der Maaten, Lawrence Chen, Liang Tan, Liz Jenkins, Louis Martin, Lovish Madaan, Lubo Malo, Lukas Blecher, Lukas Landzaat, Luke de Oliveira, Madeline Muzzi, Mahesh Pasupuleti, Mannat Singh, Manohar Paluri, Marcin Kardas, Maria Tsimpoukelli, Mathew Oldham, Mathieu Rita, Maya Pavlova, Melanie Kambadur, Mike Lewis, Min Si, Mitesh Kumar Singh, Mona Hassan, Naman Goyal, Narjes Torabi, Nikolay Bashlykov, Nikolay Bogoychev, Niladri Chatterji, Ning Zhang, Olivier Duchenne, Onur Çelebi, Patrick Alrassy, Pengchuan Zhang, Pengwei Li, Petar Vasic, Peter Weng, Prajjwal Bhargava, Pratik Dubal, Praveen Krishnan, Punit Singh Koura, Puxin Xu, Qing He, Qingxiao Dong, Ragavan Srinivasan, Raj Ganapathy, Ramon Calderer, Ricardo Silveira Cabral, Robert Stojnic, Roberta Raileanu, Rohan Maheswari, Rohit Girdhar, Rohit Patel, Romain Sauvestre, Ronnie Polidoro, Roshan Sumbaly, Ross Taylor, Ruan Silva, Rui Hou, Rui Wang, Saghar Hosseini, Sahana Chennabasappa, Sanjay Singh, Sean Bell, Seohyun Sonia Kim, Sergey Edunov, Shaoliang Nie, Sharan Narang, Sharath Raparthy, Sheng Shen, Shengye Wan, Shruti Bhosale, Shun Zhang, Simon Vandenhende, Soumya Batra, Spencer Whitman, Sten Sootla, Stephane Collot, Suchin Gururangan, Sydney Borodinsky, Tamar Herman, Tara Fowler, Tarek Sheasha, Thomas Georgiou, Thomas Scialom, Tobias Speckbacher, Todor Mihaylov, Tong Xiao, Ujjwal Karn, Vedanuj Goswami, Vibhor Gupta, Vignesh Ramanathan, Viktor Kerkez, Vincent Gonguet, Virginie Do, Vish Vogeti, Vítor Albiero, Vladan Petrovic, Weiwei Chu, Wenhan Xiong, Wenyin Fu, Whitney Meers, Xavier Martinet, Xiaodong Wang, Xiaofang Wang, Xiaoqing Ellen Tan, Xide Xia, Xinfeng Xie, Xuchao Jia, Xuewei Wang, Yaelle Goldschlag, Yashesh Gaur, Yasmine Babaei, Yi Wen, Yiwen Song, Yuchen Zhang, Yue Li, Yuning Mao, Zacharie Delpierre Coudert, Zheng Yan, Zhengxing Chen, Zoe Papakipos, Aaditya Singh, Aayushi Srivastava, Abha Jain, Adam Kelsey, Adam Shajnfeld, Adithya Gangidi, Adolfo Victoria, Ahuva Goldstand, Ajay Menon, Ajay Sharma, Alex Boesenberg, Alexei Baevski, Allie Feinstein, Amanda Kallet, Amit Sangani, Amos Teo, Anam Yunus, Andrei Lupu, Andres Alvarado, Andrew Caples, Andrew Gu, Andrew Ho, Andrew Poulton, Andrew Ryan, Ankit Ramchandani, Annie Dong, Annie Franco, Anuj Goyal, Aparajita Saraf, Arkabandhu Chowdhury, Ashley Gabriel, Ashwin Bharambe, Assaf Eisenman, Azadeh Yazdan, Beau James, Ben Maurer, Benjamin Leonhardi, Bernie Huang, Beth Loyd, Beto De Paola, Bhargavi Paranjape, Bing Liu, Bo Wu, Boyu Ni, Braden Hancock, Bram Wasti, Brandon Spence, Brani Stojkovic, Brian Gamido, Britt Montalvo, Carl Parker, Carly Burton, Catalina Mejia, Ce Liu, Changhan Wang, Changkyu Kim, Chao Zhou, Chester Hu, Ching-Hsiang Chu, Chris Cai, Chris Tindal, Christoph Feichtenhofer, Cynthia Gao, Damon Civin, Dana Beaty, Daniel Kreymer, Daniel Li, David Adkins, David Xu, Davide Testuggine, Delia David, Devi Parikh, Diana Liskovich, Didem Foss, Dingkang Wang, Duc Le, Dustin Holland, Edward Dowling, Eissa Jamil, Elaine Montgomery, Eleonora Presani, Emily Hahn, Emily Wood, Eric-Tuan Le, Erik Brinkman, Esteban Arcaute, Evan Dunbar, Evan Smothers, Fei Sun, Felix Kreuk, Feng Tian, Filippos Kokkinos, Firat Ozgenel, Francesco Caggioni, Frank Kanayet, Frank Seide, Gabriela Medina Florez, Gabriella Schwarz, Gada Badeer, Georgia Swee, Gil Halpern, Grant Herman, Grigory Sizov, Guangyi, Zhang, Guna Lakshminarayanan, Hakan Inan, Hamid Shojanazeri, Han Zou, Hannah Wang, Hanwen Zha, Haroun Habeeb, Harrison Rudolph, Helen Suk, Henry Aspegren, Hunter Goldman, Hongyuan Zhan, Ibrahim Damlaj, Igor Molybog, Igor Tufanov, Ilias Leontiadis, Irina-Elena Veliche, Itai Gat, Jake Weissman, James Geboski, James

596

597

598

600

601

602

603

604

605

606

607

608

609

610

611

612

613

614

615

616

617

618

619

620

621

622

623

624 625

626

627

628

629

630

631 632

633

634

635 636

637

638 639

640

641

642 643

644

645

646

647

Kohli, Janice Lam, Japhet Asher, Jean-Baptiste Gaya, Jeff Marcus, Jeff Tang, Jennifer Chan, Jenny Zhen, Jeremy Reizenstein, Jeremy Teboul, Jessica Zhong, Jian Jin, Jingyi Yang, Joe Cummings, Jon Carvill, Jon Shepard, Jonathan McPhie, Jonathan Torres, Josh Ginsburg, Junjie Wang, Kai Wu, Kam Hou U, Karan Saxena, Kartikay Khandelwal, Katayoun Zand, Kathy Matosich, Kaushik Veeraraghavan, Kelly Michelena, Keqian Li, Kiran Jagadeesh, Kun Huang, Kunal Chawla, Kyle Huang, Lailin Chen, Lakshya Garg, Lavender A, Leandro Silva, Lee Bell, Lei Zhang, Liangpeng Guo, Licheng Yu, Liron Moshkovich, Luca Wehrstedt, Madian Khabsa, Manav Avalani, Manish Bhatt, Martynas Mankus, Matan Hasson, Matthew Lennie, Matthias Reso, Maxim Groshev, Maxim Naumov, Maya Lathi, Meghan Keneally, Miao Liu, Michael L. Seltzer, Michal Valko, Michelle Restrepo, Mihir Patel, Mik Vyatskov, Mikayel Samvelyan, Mike Clark, Mike Macey, Mike Wang, Miquel Jubert Hermoso, Mo Metanat, Mohammad Rastegari, Munish Bansal, Nandhini Santhanam, Natascha Parks, Natasha White, Navyata Bawa, Nayan Singhal, Nick Egebo, Nicolas Usunier, Nikhil Mehta, Nikolay Pavlovich Laptev, Ning Dong, Norman Cheng, Oleg Chernoguz, Olivia Hart, Omkar Salpekar, Ozlem Kalinli, Parkin Kent, Parth Parekh, Paul Saab, Pavan Balaji, Pedro Rittner, Philip Bontrager, Pierre Roux, Piotr Dollar, Polina Zvyagina, Prashant Ratanchandani, Pritish Yuvraj, Qian Liang, Rachad Alao, Rachel Rodriguez, Rafi Ayub, Raghotham Murthy, Raghu Nayani, Rahul Mitra, Rangaprabhu Parthasarathy, Raymond Li, Rebekkah Hogan, Robin Battey, Rocky Wang, Russ Howes, Ruty Rinott, Sachin Mehta, Sachin Siby, Sai Jayesh Bondu, Samyak Datta, Sara Chugh, Sara Hunt, Sargun Dhillon, Sasha Sidorov, Satadru Pan, Saurabh Mahajan, Saurabh Verma, Seiji Yamamoto, Sharadh Ramaswamy, Shaun Lindsay, Shaun Lindsay, Sheng Feng, Shenghao Lin, Shengxin Cindy Zha, Shishir Patil, Shiva Shankar, Shuqiang Zhang, Shuqiang Zhang, Sinong Wang, Sneha Agarwal, Soji Sajuyigbe, Soumith Chintala, Stephanie Max, Stephen Chen, Steve Kehoe, Steve Satterfield, Sudarshan Govindaprasad, Sumit Gupta, Summer Deng, Sungmin Cho, Sunny Virk, Suraj Subramanian, Sy Choudhury, Sydney Goldman, Tal Remez, Tamar Glaser, Tamara Best, Thilo Koehler, Thomas Robinson, Tianhe Li, Tianjun Zhang, Tim Matthews, Timothy Chou, Tzook Shaked, Varun Vontimitta, Victoria Ajayi, Victoria Montanez, Vijai Mohan, Vinay Satish Kumar, Vishal Mangla, Vlad Ionescu, Vlad Poenaru, Vlad Tiberiu Mihailescu, Vladimir Ivanov, Wei Li, Wenchen Wang, Wenwen Jiang, Wes Bouaziz, Will Constable, Xiaocheng Tang, Xiaojian Wu, Xiaolan Wang, Xilun Wu, Xinbo Gao, Yaniv Kleinman, Yanjun Chen, Ye Hu, Ye Jia, Ye Qi, Yenda Li, Yilin Zhang, Ying Zhang, Yossi Adi, Youngjin Nam, Yu, Wang, Yu Zhao, Yuchen Hao, Yundi Qian, Yunlu Li, Yuzi He, Zach Rait, Zachary DeVito, Zef Rosnbrick, Zhaoduo Wen, Zhenyu Yang, Zhiwei Zhao, and Zhiyu Ma. The llama 3 herd of models, 2024. URL https://arxiv.org/abs/2407.21783.

- Kanav Gupta, Neha Jawalkar, Ananta Mukherjee, Nishanth Chandran, Divya Gupta, Ashish Panwar, and Rahul Sharma. SIGMA: Secure GPT inference with function secret sharing. Cryptology ePrint Archive, Paper 2023/1269, 2023. URL https://eprint.iacr.org/2023/1269.
- Meng Hao, Hongwei Li, Hanxiao Chen, Pengzhi Xing, Guowen Xu, and Tianwei Zhang. Iron: Private inference on transformers. In *Advances in Neural Information Processing Systems*, volume 35, pp. 15718–15731, 2022.
- Zhicong Huang, Wen jie Lu, Cheng Hong, and Jiansheng Ding. Cheetah: Lean and fast secure two-party deep neural network inference. In 31st USENIX Security Symposium (USENIX Security 22), pp. 809–826, 2022.
- Yuval Ishai, Eyal Kushilevitz, Rafail Ostrovsky, and Amit Sahai. Zero-knowledge from secure multiparty computation. In *Proceedings of the thirty-ninth annual ACM symposium on Theory of computing*, pp. 21–30, 2007.
- Patrick Jauernig, Ahmad-Reza Sadeghi, and Emmanuel Stapf. Trusted execution environments: Properties, applications, and challenges. *IEEE Security & Privacy*, 18(2):56–60, 2020. doi: 10.1109/MSEC.2019.2947124.
- B. Knott, S. Venkataraman, A.Y. Hannun, S. Sengupta, M. Ibrahim, and L.J.P. van der Maaten. Crypten: Secure multi-party computation meets machine learning. In *arXiv* 2109.00984, 2021.
- Dacheng Li, Rulin Shao, Hongyi Wang, Han Guo, Eric P. Xing, and Hao Zhang. Mpcformer: fast, performant and private transformer inference with mpc, 2023. URL https://arxiv.org/abs/2211.01452.

- Zhengyi Li, Kang Yang, Jin Tan, Wen jie Lu, Haoqi Wu, Xiao Wang, Yu Yu, Derun Zhao, Yancheng Zheng, Minyi Guo, and Jingwen Leng. Nimbus: Secure and efficient two-party inference for transformers, 2024. URL https://arxiv.org/abs/2411.15707.
  - Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization, 2019. URL https://arxiv.org/abs/1711.05101.
  - Anton Lozhkov, Loubna Ben Allal, Leandro von Werra, and Thomas Wolf. Fineweb-edu: the finest collection of educational content, 2024. URL https://huggingface.co/datasets/HuggingFaceFW/fineweb-edu.
  - Jinglong Luo, Guanzhong Chen, Yehong Zhang, Shiyu Liu, Hui Wang, Yue Yu, Xun Zhou, Yuan Qi, and Zenglin Xu. Centaur: Bridging the impossible trinity of privacy, efficiency, and performance in privacy-preserving transformer inference, 2024. URL https://arxiv.org/abs/2412.10652.
  - Jungho Moon, Dongwoo Yoo, Xiaoqian Jiang, and Miran Kim. THOR: Secure transformer inference with homomorphic encryption. Cryptology ePrint Archive, Paper 2024/1881, 2024. URL https://eprint.iacr.org/2024/1881.
  - Krishna Giri Narra, Zhifeng Lin, Yongqin Wang, Keshav Balasubramaniam, and Murali Annavaram. Privacy-preserving inference in machine learning services using trusted execution environments. arXiv preprint arXiv:1912.03485, 2019.
  - Jack Min Ong, Matthew Di Ferrante, Aaron Pazdera, Ryan Garner, Sami Jaghouar, Manveer Basra, Max Ryabinin, and Johannes Hagemann. Toploc: A locality sensitive hashing scheme for trustless verifiable inference, 2025. URL https://arxiv.org/abs/2501.16007.
  - Qi Pang, Jinhao Zhu, Helen Möllering, Wenting Zheng, and Thomas Schneider. BOLT: Privacy-preserving, accurate and efficient inference for transformers. Cryptology ePrint Archive, Paper 2023/1893, 2023. URL https://eprint.iacr.org/2023/1893.
  - Wenjie Qu, Yijun Sun, Xuanming Liu, Tao Lu, Yanpei Guo, Kai Chen, and Jiaheng Zhang. zkgpt: An efficient non-interactive zero-knowledge proof framework for llm inference. In *34st USENIX Security Symposium (USENIX Security 25)*, 2025.
  - Qwen, :, An Yang, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chengyuan Li, Dayiheng Liu, Fei Huang, Haoran Wei, Huan Lin, Jian Yang, Jianhong Tu, Jianwei Zhang, Jianxin Yang, Jiaxi Yang, Jingren Zhou, Junyang Lin, Kai Dang, Keming Lu, Keqin Bao, Kexin Yang, Le Yu, Mei Li, Mingfeng Xue, Pei Zhang, Qin Zhu, Rui Men, Runji Lin, Tianhao Li, Tianyi Tang, Tingyu Xia, Xingzhang Ren, Xuancheng Ren, Yang Fan, Yang Su, Yichang Zhang, Yu Wan, Yuqiong Liu, Zeyu Cui, Zhenru Zhang, and Zihan Qiu. Qwen2.5 technical report, 2025. URL https://arxiv.org/abs/2412.15115.
  - Mohamed Sabt, Mohammed Achemlal, and Abdelmadjid Bouabdallah. Trusted execution environment: What it is, and what it is not. In 2015 IEEE Trustcom/BigDataSE/Ispa, volume 1, pp. 57–64. IEEE, 2015.
  - Sanjif Shanmugavelu, Mathieu Taillefumier, Christopher Culver, Oscar Hernandez, Mark Coletti, and Ada Sedova. Impacts of floating-point non-associativity on reproducibility for hpc and deep learning applications, 2024. URL https://arxiv.org/abs/2408.05148.
  - Haochen Sun, Jason Li, and Hongyang Zhang. zkllm: Zero knowledge proofs for large language models, 2024. URL https://arxiv.org/abs/2404.16109.
  - Yifan Sun, Yuhang Li, Yue Zhang, Yuchen Jin, and Huan Zhang. Svip: Towards verifiable inference of open-source large language models, 2025. URL https://arxiv.org/abs/2410.22307.
  - Kimi Team, Yifan Bai, Yiping Bao, Guanduo Chen, Jiahao Chen, Ningxin Chen, Ruijue Chen, Yanru Chen, Yuankun Chen, Yutian Chen, Zhuofu Chen, Jialei Cui, Hao Ding, Mengnan Dong, Angang Du, Chenzhuang Du, Dikang Du, Yulun Du, Yu Fan, Yichen Feng, Kelin Fu, Bofei Gao, Hongcheng Gao, Peizhong Gao, Tong Gao, Xinran Gu, Longyu Guan, Haiqing Guo, Jianhang Guo, Hao Hu, Xiaoru Hao, Tianhong He, Weiran He, Wenyang He, Chao Hong, Yangyang Hu,

Zhenxing Hu, Weixiao Huang, Zhiqi Huang, Zihao Huang, Tao Jiang, Zhejun Jiang, Xinyi Jin, Yongsheng Kang, Guokun Lai, Cheng Li, Fang Li, Haoyang Li, Ming Li, Wentao Li, Yanhao Li, Yiwei Li, Zhaowei Li, Zheming Li, Hongzhan Lin, Xiaohan Lin, Zongyu Lin, Chengyin Liu, Chenyu Liu, Hongzhang Liu, Jingyuan Liu, Junqi Liu, Liang Liu, Shaowei Liu, T. Y. Liu, Tianwei Liu, Weizhou Liu, Yangyang Liu, Yibo Liu, Yiping Liu, Yue Liu, Zhengying Liu, Enzhe Lu, Lijun Lu, Shengling Ma, Xinyu Ma, Yingwei Ma, Shaoguang Mao, Jie Mei, Xin Men, Yibo Miao, Siyuan Pan, Yebo Peng, Ruoyu Qin, Bowen Qu, Zeyu Shang, Lidong Shi, Shengyuan Shi, Feifan Song, Jianlin Su, Zhengyuan Su, Xinjie Sun, Flood Sung, Heyi Tang, Jiawen Tao, Qifeng Teng, Chensi Wang, Dinglu Wang, Feng Wang, Haiming Wang, Jianzhou Wang, Jiaxing Wang, Jinhong Wang, Shengjie Wang, Shuyi Wang, Yao Wang, Yejie Wang, Yiqin Wang, Yuxin Wang, Yuzhi Wang, Zhaoji Wang, Zhengtao Wang, Zhexu Wang, Chu Wei, Qianqian Wei, Wenhao Wu, Xingzhe Wu, Yuxin Wu, Chenjun Xiao, Xiaotong Xie, Weimin Xiong, Boyu Xu, Jing Xu, Jinjing Xu, L. H. Xu, Lin Xu, Suting Xu, Weixin Xu, Xinran Xu, Yangchuan Xu, Ziyao Xu, Junjie Yan, Yuzi Yan, Xiaofei Yang, Ying Yang, Zhen Yang, Zhilin Yang, Zonghan Yang, Haotian Yao, Xingcheng Yao, Wenjie Ye, Zhuorui Ye, Bohong Yin, Longhui Yu, Enming Yuan, Hongbang Yuan, Mengjie Yuan, Haobing Zhan, Dehao Zhang, Hao Zhang, Wanlu Zhang, Xiaobin Zhang, Yangkun Zhang, Yizhi Zhang, Yongting Zhang, Yu Zhang, Yutao Zhang, Yutong Zhang, Zheng Zhang, Haotian Zhao, Yikai Zhao, Huabin Zheng, Shaojie Zheng, Jianren Zhou, Xinyu Zhou, Zaida Zhou, Zhen Zhu, Weiyu Zhuang, and Xinxing Zu. Kimi k2: Open agentic intelligence, 2025. URL https://arxiv.org/abs/2507.20534.

Rahul Thomas, Louai Zahran, Erica Choi, Akilesh Potti, Micah Goldblum, and Arka Pal. Cascade: Token-sharded private llm inference, 2025. URL https://arxiv.org/abs/2507.05228.

Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, Dan Bikel, Lukas Blecher, Cristian Canton Ferrer, Moya Chen, Guillem Cucurull, David Esiobu, Jude Fernandes, Jeremy Fu, Wenyin Fu, Brian Fuller, Cynthia Gao, Vedanuj Goswami, Naman Goyal, Anthony Hartshorn, Saghar Hosseini, Rui Hou, Hakan Inan, Marcin Kardas, Viktor Kerkez, Madian Khabsa, Isabel Kloumann, Artem Korenev, Punit Singh Koura, Marie-Anne Lachaux, Thibaut Lavril, Jenya Lee, Diana Liskovich, Yinghai Lu, Yuning Mao, Xavier Martinet, Todor Mihaylov, Pushkar Mishra, Igor Molybog, Yixin Nie, Andrew Poulton, Jeremy Reizenstein, Rashi Rungta, Kalyan Saladi, Alan Schelten, Ruan Silva, Eric Michael Smith, Ranjan Subramanian, Xiaoqing Ellen Tan, Binh Tang, Ross Taylor, Adina Williams, Jian Xiang Kuan, Puxin Xu, Zheng Yan, Iliyan Zarov, Yuchen Zhang, Angela Fan, Melanie Kambadur, Sharan Narang, Aurelien Rodriguez, Robert Stojnic, Sergey Edunov, and Thomas Scialom. Llama 2: Open foundation and fine-tuned chat models, 2023. URL https://arxiv.org/abs/2307.09288.

Colin White, Samuel Dooley, Manley Roberts, Arka Pal, Ben Feuer, Siddhartha Jain, Ravid Shwartz-Ziv, Neel Jain, Khalid Saifullah, Sreemanti Dey, Shubh-Agrawal, Sandeep Singh Sandha, Siddartha Naidu, Chinmay Hegde, Yann LeCun, Tom Goldstein, Willie Neiswanger, and Micah Goldblum. Livebench: A challenging, contamination-limited llm benchmark, 2025. URL https://arxiv.org/abs/2406.19314.

Andrew C. Yao. Protocols for secure computations. In 23rd Annual Symposium on Foundations of Computer Science (sfcs 1982), pp. 160–164, 1982. doi: 10.1109/SFCS.1982.38.

Mu Yuan, Lan Zhang, and Xiang-Yang Li. Secure transformer inference protocol, 2024. URL https://arxiv.org/abs/2312.00025.

Jiawen Zhang, Xinpeng Yang, Lipeng He, Kejia Chen, Wen jie Lu, Yinghao Wang, Xiaoyang Hou, Jian Liu, Kui Ren, and Xiaohu Yang. Secure transformer inference made non-interactive. Cryptology ePrint Archive, Paper 2024/136, 2024. URL https://eprint.iacr.org/2024/136.

Fei Zheng, Chaochao Chen, Zhongxuan Han, and Xiaolin Zheng. Permllm: Private inference of large language models within 3 seconds under wan, 2024. URL https://arxiv.org/abs/2405.18744.

# A BACKGROUND AND RELATED WORK

In this section, we provide a brief background on general methods of privacy-preserving function computation, general methods of verification, and their application to LLM inference in particular.

#### A.1 PRIVACY-PRESERVATION

There are four main families of privacy-preserving inference of LLMs that have been proposed in the literature: **SMPC** (Secure Multi-Party Computation), **FHE** (Fully Homomorphic Encryption), **TEEs** (Trusted Execution Environments), and **statistical methods**. Here we provide brief background on each of these.

SMPC SMPC protocols split the required computation among multiple parties. The key ideas were originally developed in the 1980s (Yao, 1982; Goldreich et al., 1987) and provide mathematical guarantees that no single party can reconstruct the data on their own. Recently, the methodologies of SMPC have been applied to LLMs (Huang et al., 2022; Hao et al., 2022; Pang et al., 2023; Akimoto et al., 2023; Dong et al., 2023; Li et al., 2024). A difficulty uniformly faced by these protocols is efficient computation of the many non-linearities present in transformer-based LLMs; most of the works attempt to ameliorate this by using piecewise polynomial approximations which are more well-suited for MPC algorithms. However, this approximation leads to degraded inference results, and remains more expensive than direct computation of the non-linearities. The requirement of multiple parties also engenders significant communication overheads, and the further non-collusion requirement among the parties may be difficult to guarantee.

FHE FHE protocols require only a single party and make use of cryptographic methods to ensure that the result of the computation on the ciphertext is the same as that performed on the plaintext. The adjective 'fully' indicates the capability of performing arbitrary computations, not limited to a particular type or complexity. The first plausible construction of an FHE scheme was described in Gentry (2009); a more modern and widely used incarnation is CKKS (Cheon et al., 2017). Recently, CKKS has been further optimized and applied to LLM inference (Moon et al., 2024; Zhang et al., 2024), but similar issues arise with the non-linearities as SMPC methods. The overheads both for linear and non-linear operations are typically even larger than those in the SMPC setting.

TEEs Trusted Execution Environments (TEEs) (Sabt et al., 2015; Narra et al., 2019) create secure and isolated enclaves at the hardware level. This ensures confidentiality via memory encryption – allowing only the process running in the enclave to read the data. Furthermore, TEEs support integrity via attestation mechanisms. However, a significant concern is the vulnerability to side-channel attacks (Jauernig et al., 2020). Furthermore, attestation is only provided at boot-time and is not equivalent to an ongoing verification process. This process typically involves the TEE measuring the code and its environment, signing these measurements cryptographically, and sending a report for external verification. However, this is often a one-time check at the start and does not guarantee the integrity of the TEE throughout its execution. Finally, in cloud environments, attestation can rely on the cloud provider's services, which means users must trust the provider's proprietary attestation process without full transparency. This introduces a level of trust in the cloud provider's integrity, as these attestation services can be opaque "black boxes" that are not open to external audit. Moreover, there may be no independent way to verify the boot measurements provided by the cloud provider's infrastructure.

**Statistical Methods** A more broad and diverse grouping than the above is what we term 'statistical methods'. These are protocols without the mathematical guarantees of FHE or SMPC approaches, or the hardware-based guarantees of TEEs, but that instead employ statistical or empirical arguments to support the difficult of reversing ciphertext. Some ideas in this domain include the use of permutation-based security (Zheng et al., 2024; Yuan et al., 2024; Luo et al., 2024) or token-sharding based security (Thomas et al., 2025). These methods typically trade off the stronger guarantees of the above methods for greatly reduced overheads, sometimes approaching similar speeds to vanilla inference.

#### A.2 VERIFICATION

**Zero-Knowledge Proofs (ZKP)** ZKPs are a class of methods that allows one party (the prover) to prove to another party (the verifier) that a statement is true, without revealing any additional

information beyond the proof itself. The main properties that ZKPs satisfy are completeness (an honest prover can convince a verifier that they performed the work as stated), soundness (a dishonest prover cannot convince a verifier that they performed the work), and the zero-knowledge property of not revealing any further information than the fact the work was done as stated. The first ZK protocol was introduced in 1985 in Goldwasser et al. (1985). Recently, ZK methods have been applied as proofs of inference for machine learning models, and specifically LLMs, in works such as Sun et al. (2024); Qu et al. (2025). However, these approaches remain thousands of times slower than vanilla inference – for example, zkLLM takes 15 minutes for generating a proof of a single forward pass for Llama-2-13B, compared to milliseconds for vanilla inference.

Statistical Methods Analogously to statistical methods of privacy-preservation, very recent work has investigated methods of relaxing the standard of proof of work provided in order to reduce computational overhead. Ong et al. (2025) encodes and validates the most salient features of the last hidden state tensor of an LLM using a compact, verifiable proof, which is then recomputed in parallel by the verifier. Although the authors demonstrate how to set up a commitment scheme that has relatively little overhead to the prover, and verification is faster than full recomputation thanks to parallelization, there is still a requirement for the verifier to perform a full LLM forward pass, potentially necessitating specialized hardware. Sun et al. (2025) proposes the use of a 'proxy task' based on the last hidden layer features of an LLM that can then be utilized by the user to compare to a label that they would expect based on their original input. The method proposed requires trust assumptions from the platform for generation of the proxy-task feature extractor and labeller networks, as well as secret generation/embedding, and adds the overhead of computation to perform all of the above.

## B FORMAL ALGORITHM FOR PROTOCOL 1

The procedure is comprised of three components: cache generation through Algorithm 1, inference request through Algorithm 2, and the verification stage through Algorithm 3.

#### **Algorithm 1** Cache Generation

```
Input: model, cache size |C| \in \mathbb{N}, sentinel token count K \in \mathbb{N}
Output: cache: mapping s \mapsto \ell_{1:K} \in \mathbb{R}^{K \times V}
 1: cache \leftarrow \emptyset
 2: while |cache| < |C| do
           s_{1:K} \leftarrow \text{sample with replacement } K \text{ tokens from } V
 3:
                                                                                             \triangleright initialize K \times K attention mask
 4:
           m_{1:K,1:K} \leftarrow 0
 5:
           for i = 1..K do
 6:
                for j = 1..i do
 7:
                      m_{i,j} \leftarrow 1
                 end for
 8:
 9:
           end for
                                                                                                                           \triangleright\,\ell \in \mathbb{R}^{K\times V}
10:
           \ell \leftarrow \text{model.forward}(s, m)
11:
           cache[s] \leftarrow \ell
12: end while
13: return cache
```

# Algorithm 2 Inference Request

864

871

872

873

874

875

876

877

878

879

880

881

882

883

884

885

886

900

901

902

903

904

905

906

907

908

909

910

911

```
Input: prompt token embeddings x_{1:N}, attention mask a_{1:N}, sentinel token sequence s_{1:K}
Output: logits \ell_{1:N+K} \in \mathbb{R}^{(N+K)\times V}
 1: positions p_{1:K} \leftarrow sample without replacement K times from Uniform[1, N + K]
 2: augmented embeddings x'_{1:N+K} \leftarrow insert sentinel tokens s_{1:K} at positions p_{1:K}
 3: augmented mask a'_{1:N+K} \leftarrow \text{expand } a_{1:N} \text{ at positions } p_{1:K} \text{ with 0-filled rows and columns}
 4: for i = 1..K do
 5:
         for j = 1..i do
              a'[p_i, p_j] \leftarrow 1
 6:
 7:
         end for
 8:
         for j=i..K do
 9:
              a'[p_j, p_i] \leftarrow 1
10:
         end for
11: end for
12: x' \leftarrow \text{encrypt}(x')
13: a' \leftarrow \text{encrypt}(a')
14: encrypted \ell_{1:N+K} \leftarrow inference provider forward pass on encrypted x', a'
15: return decrypt(\ell_{1:N+K})
```

# Algorithm 3 Verification

```
Input: logits \ell_{1:N+K} \in \mathbb{R}^{(N+K)\times V}, sentinel positions p_{1:K} \subset \{1, 2, ...N + K\}, sentinel sequence
     s_{1:K}, cache \in \mathbb{R}^{C \times K \times V}, tolerance tol > 0
Output: verified: bool
 1: verified \leftarrow true
 2: for i = 1..K do
 3:
          p' \leftarrow p[i]
          err \leftarrow \|\ell[p'] - \operatorname{cache}[s][i]\|_1
 4:
          if err > tol then
 5:
                verified \leftarrow false
 6:
 7:
          end if
 8: end for
 9: return verified
```

# C FORMAL ALGORITHM FOR PROTOCOL 2

The procedure is also comprised of three algorithms similar to those in Section B: noised embedding generation Algorithm 4, noisy inference request Algorithm 5, and verification with noisy prediction Algorithm 6.

## Algorithm 4 Noised Embedding Generation

918

919 920

921

922

923 924

925

926

927

928

929

930

931932933

934

935

936

937

938

939

940

941

942

943

944

945

946

947

948

949

951

952

953

954

955

970 971

```
Input: discrete noise set B, trained NoiseEmbedder E, embedding d_e, token embedding e \in R^{d_e}

Output: sampled noise b \in R^{d_e}, noised embedding e_b \in \mathbb{R}^{d_e}

1: b \leftarrow sample one value uniformly from B

2: b_e \leftarrow \operatorname{concat}(e, b) \triangleright b_e \in \mathbb{R}^{2d_e}

3: e' \leftarrow E.forward(b_e) \triangleright e' \in \mathbb{R}^{d_e}

4: return (b, e')
```

### **Algorithm 5** Noisy Inference Request

```
Input: prompt token embeddings x_{1:N}, attention mask a_{1:N}, sentinel token sequence s_{1:K}
Output: hidden states h \in \mathbb{R}^{(N+K)\times d_h}, noise_cache \in B^N
 1: positions p_{1:K} \leftarrow sample without replacement K times from Uniform[1, N + K]
 2: augmented embeddings x'_{1:N+K} \leftarrow insert sentinel tokens s_{1:K} at positions p_{1:K}
 3: augmented mask a'_{1:N+K} \leftarrow \text{expand } a_{1:N} at positions p_{1:K} with 0-filled rows and columns
 4: for i = 1...K do
         for j = 1..i do
 5:
             a'[p_i, p_j] \leftarrow 1
 6:
 7:
         end for
 8:
         for j=i..K do
 9:
              a'[p_i, p_i] \leftarrow 1
10:
         end for
11: end for
12: for each non-sentinel token t in x' do
         b, e' \leftarrow \text{call Algorithm 4 on } x'[t]
13:
         x'[t] \leftarrow e'
14:
15:
         noise\_cache[t] \leftarrow b
16: end for
17: x' \leftarrow \text{encrypt}(x')
18: a' \leftarrow \text{encrypt}(a')
19: encrypted h_{1:N+K} \leftarrow inference provider forward pass on encrypted x', a'
20: h \leftarrow decrypt(h_{1:N+K})
21: return h, noise_cache
```

# Algorithm 6 Verification With Noise Prediction

972

973

974 975

976

977

978

979

980

981

982

983

984

985

986

987

990

991

992

993

994995996

997 998

999 1000

1001

1002

1003

1005

1007

1008

1009

1010 1011 1012

1013 1014

1015 1016

1017

1018

1019 1020

1021

1023

1024

1025

```
Input: decrypted hidden states h_{1:N+K} \in \mathbb{R}^{(N+K)\times d_h}, sentinel positions p_{1:K}, sentinel sequence s_{1:K}, logit_cache \in R^{C\times K\times V}, noise_cache \in B^N, NoisePredictor NP, logit projection L:
      \mathbb{R}^{d_h} \to \mathbb{R}^V, sentinel tolerance tol_s
Output: verified: bool
 1: for i = 1..K do
           p' \leftarrow p_i
 2:
 3:
           \ell \leftarrow L(h[p'])
 4:
           err \leftarrow \|\ell - \text{logit\_cache}[s][i]\|_1
 5:
           if err > tol_s then
                 verified \leftarrow false
 6:
 7:
           end if
 8: end for
 9: for j = 1..N + K do
           if j \in P then
10:
                 continue
11:
           end if
12:
           b \leftarrow NP.\text{forward}(h[j])
13:
14:
           if b \neq \text{noise\_cache}[j] then
15:
                  verified \leftarrow false
16:
           end if
17: end for
18: return verified
```

# D Training Details for Protocol 2

In this section we provide further details for the experiments conducted in Section 5.3.

We train on the FineWeb-Edu dataset (Lozhkov et al., 2024). This is a large-scale dataset of 1.3T total tokens consisting of high-quality educational web pages filtered from the larger FineWeb dataset. This dataset has been used for pretraining, and is suitable for general testing of language modeling capabilities of LLMs. We take 40000 samples from this dataset and divide these into an 80/20 train-validation split. We perform training for 500 steps with a batch size of 64 on sequences of length 256. We use the optimizer AdamW (Loshchilov & Hutter, 2019) with a learning rate of 5e-4, with no warmup steps. The base model is Llama-3.2-1B, and the weights of this model are frozen; gradients are backpropagated through this model in order to reach the NoiseEmbedder module.

We utilize the same sampled noise at every token position, but a different noise is sampled for each batch element. Our results are reported using a  $\lambda$  hyperparameter value of 3.5. We train using an A100 GPU.

#### E PROTOCOL 3: KEY APPENDING – FURTHER DETAILS

In this section, we give further details on the Key Appending protocol introduced in Section 6.

We also tested the same protocol but instead of appending the key repeating prompt to the end of the user's prompt, we insert it in any whitespace randomly. This approach performs less well – we report results for this in Section F.

# E.1 FURTHER PROTOCOL DETAILS

In addition to the augmented prompt, we also modify the system prompt to encourage model compliance to the protocol.

**System prompt.** The user adds a system-level instruction that enforces the verification protocol and prevents premature conversation termination:

Table 3: Transcription success rate on 1000 prompts of our 'key appending' verification protocol, with keys of length K=3 words. We see that models with parameter sizes of 3B and above obtain very high transcription rates of > 98%.

Model	Transcription rate
Llama 3.2 1B	56.6%
Gemma 3 1B	73.1%
Llama 3.2 3B	98.1%
Gemma 3 4B	99.7%
Mistral 7B	99.6%
Llama 3.1 8B	98.7%
Gemma 3 12B	98.0%
Mistral 24B	100.0%
Qwen 2.5 32B	99.9%
Llama 3.1 70B	99.6%

Table 4: LiveBench scores for models under the 'key appending' protocol. Higher is better. There is a relatively large degradation in performance for smaller models, but models of size 8B and larger exhibit much smaller relative degradation.

Model	Vanilla	Append	$\Delta$
Llama 3.2 1B	10.7	6.7	-4.0
Gemma 3 1B	14.7	9.9	-4.8
Llama 3.2 3B	20.5	16.9	-3.6
Qwen 2.5 3B	24.2	18.3	-5.9
Gemma 3 4B	30.2	26.4	-3.8
Mistral 7B	20.4	14.5	-5.9
Llama 3.1 8B	25.4	25.6	+0.2
Gemma 3 12B	41.0	36.8	-4.2
Mistral 24B	30.5	32.1	+1.6
Qwen 2.5 32B	42.7	41.7	-1.0
Llama 3.1 70B	42.3	39.8	-2.5

You are a helpful assistant who should never speak in two consecutive turns. At the end of your response, repeat the key mentioned at the end of the prompt. You must print the key between tags like the following structure: <key> \*insert key here\* </key>.

This system prompt explicitly conditions the model to always conclude with the verification key enclosed in  $\langle \text{key} \rangle$  . . .  $\langle \text{key} \rangle$  tags, and discourages it from generating the turn-termination token prior to key emission.

**Stopping criterion.** Unlike standard decoding, where inference continues until the model produces an end-of-sequence token, the user adopts a custom stopping rule: decoding halts once the closing tag </key> is generated. This ensures that inference completes exactly after the verification key is produced, with no trailing tokens.

#### E.2 EXPERIMENTS

**Key Transcription Capability.** We first examine the capability of LLMs to perform our protocol successfully; this is analogous to the cryptographic property of *completeness*. We evaluate multiple open-source models on a random sample of 1000 prompts from the Databricks dolly-15k dataset (Argilla, 2023). The key is sampled uniformly at random from the standard Ubuntu words file provided by the wordlist package (any similar list of English words suffices) and appended via the protocol; success is recorded if the words enclosed in the key tags match the key exactly.

Table 3 reports the transcription success rate. We observe near-perfect adherence for models at or above the 3–4B scale. Therefore, a simple capacity criterion suffices in practice: models with  $\geq$ 3B parameters reliably satisfy the protocol's instruction, making them suitable drop-in choices for verified inference with key appending.

**Downstream Performance Impact.** We now test the impact of the additional key-transcription instruction on model downstream performance. To quantify this, we evaluate all models on LIVEBENCH (White et al., 2025) (30–05–2025 release), a benchmark testing model performance on a range of tasks including data analysis, instruction following, language, math, and reasoning. Table 4 compares overall performance of the models in vanilla inference against our protocol. Extended tables showing complete results for all LiveBench categories are presented in Section H and Section F.

We find that there there is a relatively large performance impact for smaller models. However, for larger models, the performance impact is reduced. Indeed, for Mistral 24B, the performance is actually slightly higher – which we attribute to the natural variability inherent in the benchmark. In

practice, using a model of size 8B or larger seems sufficient to ensure minimal relative downstream performance impact from applying this protocol.

**Approximation Attack Experiments** As described in Section 6.2, this protocol is potentially vulnerable to model approximation attacks. We perform an in-depth examination of the viability of such in the specific case of use of SMPC for privacy preservation, with one honest party, in Section G. We find that in this setting, any such approximation does fail.

### F RANDOM WHITESPACE PROTOCOL

As an alternative to the key appending protocol of Section E, we also tested requiring the model to repeat the key in a random whitespace. **Protocol description.** This protocol adopts the same stopping criterion and verification method as the key–appending protocol, but modifies both the system prompt and the user prompt transformation.

SYSTEM PROMPT. We prepend the following system instruction:

You are a helpful assistant who should never speak in two consecutive turns. Somewhere in the prompt, you will find a secret key of {K} words in curly braces right after the marker <{marker}>. Print that key at the very end of your response. Ignore what is inside the <verify> </verify> tags when answering the main prompt.\nYou must print the key between tags like the following structure: <key> \*insert key here\* </key>.

This system prompt enforces the placement of the key at the end of the model's response and ensures that the verification tags are ignored during the main task, preventing interference with downstream output.

USER PROMPT TRANSFORMATION. Given an original user prompt p, the user selects a random whitespace location and inserts the following structure:

```
<verify> <{marker}> {key} </{marker}> </verify>.
```

Here, the marker is a randomly generated four-character ASCII string, and the key consists of three English words sampled uniformly at random, as in the key-appending protocol.

DESIGN RATIONALE. The system prompt explicitly instructs the model to ignore the inserted tags when answering the main query, which tries to minimize the impact of the injected verification key on downstream task performance. Moreover, we deliberately employ HTML-like tags for three reasons:

- 1. Large language models are extensively exposed during pretraining to HTML/XML patterns, which aids reliable parsing and generation.
- 2. Wrapping the marker–key pair inside <verify> tags avoids accidental collisions with ordinary prompts (e.g., programming queries that might already include custom markers).
- 3. Randomly generating the marker string reduces the probability of unintentional matches with existing content, while including the outer <verify> tags improves transcription accuracy compared to using only <marker>... </marker>... </marker>...

The transcription rates and downstream performance impact of the protocol are shown in Table 5 and Table 6 respectively. Although the transcription rates match that of key appending for models of size 8B and above, the downstream performance impact is significantly larger than that of key appending.

Table 6: LiveBench scores for models under the 'random whitespace' protocol. Higher is better. There is generally a larger degradation in performance than for the key appending protocol.

Model	Vanilla	Random Whitespace	$\Delta$
Llama 3.2 1B	10.7	5.4	-5.3
Gemma 3 1B	14.7	9.2	-5.5
Llama 3.2 3B	20.5	14.2	-6.3
Qwen 2.5 3B	24.2	17.2	-7.0
Gemma 3 4B	30.2	21.3	-8.9
Mistral 7B	20.4	12.2	-8.2
Llama 3.1 8B	25.4	20.7	-4.7
Gemma 3 12B	41.0	29.1	-11.9
Mistral 24B	30.5	25.3	-5.2
Qwen 2.5 32B	42.7	39.3	-3.4
Llama 3.1 70B	42.3	32.1	-10.2

Table 5: Transcription success rate on 1000 prompts of our 'Random Whitespace' verification protocol, with keys of length K=3 words. We see that models with parameter sizes of 8B and above obtain very high transcription rates of >98%.

Model	Transcription rate
Llama 3.2 1B	6.7%
Gemma 3 1B	2.2%
Llama 3.2 3B	88.8%
Gemma 3 4B	79.0%
Mistral 7B	86.6%
Llama 3.1 8B	98.5%
Gemma 3 12B	99.0%
Mistral 24B	99.6%
Qwen 2.5 32B	99.2%
Llama 3.1 70B	99.3%

#### G APPROXIMATION ATTACK EXPERIMENTS – KEY APPENDING

We perform approximation attack tests in the SMPC setting. We assume the existence of at least one honest party; in the case where all parties are dishonest (i.e. performing the same, matching approximation), the approximated model still can potentially accurately produce the key and evade detection. We use the CrypTen Python library (Knott et al., 2021).

For the dishonest party, we uniformly reduce the rank of all weight matrices in the models to various proportions of the original rank, and test the protocol to see whether the approximated model is still capable of correctly outputting the key. We select the following for our parameters:

- 1. Models: Llama 3.2 3B Instruct, Qwen 2.5 3B. We select two models with very high key transcription rates in the non-attack setting from different model classes.
- 2. We test the reduction of original ranks of very weight matrix M to the following percentage reductions: 1%, 5%, 10%, 25%, 50%, 75%, 90%, 99%, using SVD. We desire to test a wide variety of different ranks, ranging from an extremely significant reduction in rank to a slight decrease in rank, and we hence select the previously listed percentages for significant coverage of all of these possibilities.

For each combination of model and rank, we run the framework as described at the beginning of this section, selecting n=20 prompts and K=3 words.

The results of such experiments revealed that regardless of the model used or rank approximated to, the model was *always unable to output the key* (i.e. 0 of the 20 tests succeeded). Notably, even in the

99% test, both models were unable to produce anything legible, and tokens outputted were entirely random: an example decoded result from one prompt was "deesestiftigiongh" with random unicode characters inserted inside.

**Quantizer Attacks** A malicious actor can also potentially quantize the model's weights to a different precision, which is straightforward to test: given a model, we quantize all its weights to a different precision and perform the common tests to determine performance.

We again perform tests in an SMPC setting encrypted with CrypTen with two parties, one honest and one dishonest. We again note the potential weakness of this strategy when both parties are dishonest or a different encryption scheme is used.

- 1. Models: Llama 3.2 3B Instruct, Qwen 2.5 3B. We select the same models as for the low-rank approximations, due to their ordinarily high transcription rates.
- 2. Precisions: 8-bit and 4-bit floats. The weights in the Llama model tested are 16-bit floats at full precision, and in the Qwen model are 32-bit floats. Therefore, to reduce precision, we test quantization to 8- and 4-bit precision.

Once again, we run the common testing framework with n=20 prompts and K=3 words. Similar to the low-rank tests, in all cases, the models were *never able to output* the key, or in fact anything legible, revealing the effectiveness of the key appending protocol in defending against quantization attacks.

### H KEY APPENDING – EXTENDED LIVEBENCH RESULTS

Table 7: LiveBench category scores for vanilla inference.

Model	Average	Data Analysis	Instr. Follow.	Language	Math	Reasoning
Llama 3.2 1B	10.7	12.9	25.1	0.0	9.5	5.8
Gemma 3 1B	14.7	12.0	42.1	3.7	13.1	2.9
Llama 3.2 3B	20.5	23.3	48.4	3.5	15.5	11.9
Qwen 2.5 3B	24.2	29.0	43.2	10.7	23.5	14.6
Gemma 3 4B	30.2	38.3	61.5	6.3	33.0	11.8
Mistral 7B	20.4	26.4	46.2	1.5	13.4	14.4
Llama 3.1 8B	25.4	36.0	48.0	13.8	15.9	13.1
Gemma 3 12B	41.0	46.4	71.2	19.3	39.5	28.8
Mistral 24B	30.5	42.1	50.4	17.3	19.0	23.4
Qwen 2.5 32B	42.7	50.7	61.2	27.3	43.9	30.4
Llama 3.1 70B	42.3	52.6	65.9	30.3	31.4	31.4

Table 8: LiveBench category scores under the 'key appending' protocol. Higher is better.

		•		<u> </u>		
Model	Average	Data Analysis	Instr. Follow.	Language	Math	Reasoning
Llama 3.2 1B	6.7	0.9	24.9	0.0	2.3	5.5
Gemma 3 1B	9.9	3.3	32.1	2.3	4.8	6.6
Llama 3.2 3B	16.9	8.0	43.4	7.8	11.7	13.5
Qwen 2.5 3B	18.3	18.2	30.4	5.8	21.4	15.9
Gemma 3 4B	26.4	38.6	41.3	8.0	24.5	19.8
Mistral 7B	14.5	21.7	31.8	6.7	6.6	5.5
Llama 3.1 8B	25.6	35.3	51.7	9.3	15.2	16.6
Gemma 3 12B	36.8	46.1	60.5	16.5	37.0	24.0
Mistral 24B	32.1	41.1	47.2	24.0	21.0	26.9
Qwen 2.5 32B	41.7	47.2	58.2	24.0	44.2	35.0
Llama 3.1 70B	39.8	50.4	69.2	22.0	29.1	28.5