

VocabTailor: Dynamic Vocabulary Selection for Downstream Tasks in Small Language Models

Hanling Zhang^{1,2*}, Yayu Zhou^{7*}, Tongcheng Fang³, Zhihang Yuan^{3†},
Guohao Dai^{2,5}, Wanli Ouyang^{1,4,6}, Yu Wang^{3‡}

¹The Chinese University of Hong Kong, ²Infinigence AI, ³Tsinghua University,

⁴SLAI, ⁵Shanghai Jiao Tong University,

⁶Shanghai Artificial Intelligence Laboratory, ⁷Independent Researcher

Abstract

Small Language Models (SLMs) provide computational advantages in resource-constrained environments, yet memory limitations remain a critical bottleneck for edge device deployment. A substantial portion of SLMs’ memory footprint stems from vocabulary-related components, particularly embeddings and language modeling (LM) heads, due to large vocabulary sizes. Existing static vocabulary pruning, while reducing memory usage, suffers from rigid, one-size-fits-all designs that cause information loss during the prefill stage and lack flexibility. In this work, we identify two key principles underlying the vocabulary reduction challenge: the *lexical locality* principle, the observation that only a small subset of tokens is required during any single inference, and the *asymmetry in computational characteristics* between vocabulary-related components of SLM. Based on these insights, we introduce **VocabTailor**, a novel decoupled dynamic vocabulary selection framework that addresses memory constraints through offloading embedding and implements a hybrid static-dynamic vocabulary selection strategy for LM Head, enabling on-demand loading of vocabulary components. Comprehensive experiments across diverse downstream tasks demonstrate that VocabTailor achieves a reduction of up to 99% in the memory usage of vocabulary-related components with minimal or no degradation in task performance, substantially outperforming existing static vocabulary pruning. Our code is available at <https://github.com/AwakenedInsects/VocabTailor>.

1 Introduction

Large Language Models (LLMs) (Brown et al., 2020; Touvron et al., 2023a,b; Achiam et al., 2023; Team et al., 2024; Anthropic, 2024, 2025; Bai et al.,

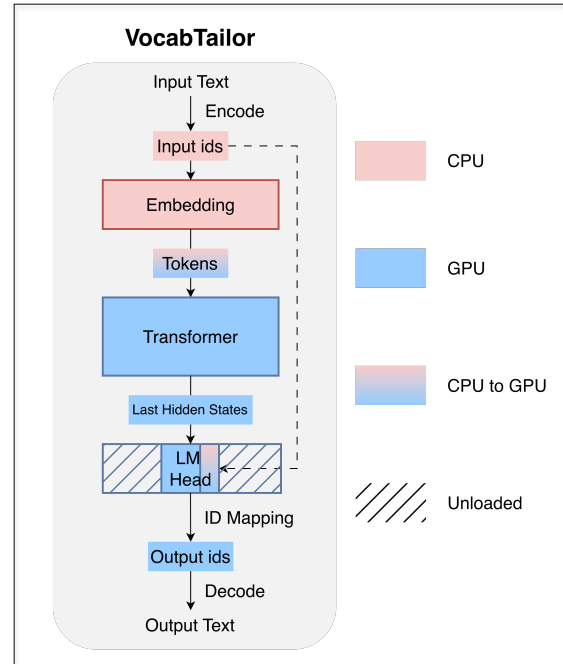


Figure 1: Overview of VocabTailor

2023; Guo et al., 2025) have rapidly become foundational to modern AI applications. Recently, increasing attention has turned towards small language models (SLMs), which are better suited for deployment on edge devices and in resource-constrained environments. Despite their compact size, memory still remains a bottleneck, particularly for edge devices with limited GPU memory. A key driver of this bottleneck is the model’s vocabulary size, which directly impacts the memory footprint of both the embedding layer and the language modeling (LM) head. For example, in the Llama 3.2 1B model with a 128K-token vocabulary, the embedding and LM head account for over 20% of the total memory usage. As SLMs are scaled down and deployed under tighter memory constraints, vocabulary-related memory inefficiencies become increasingly unsustainable, posing a

*Equal contribution.

†Co-corresponding author.

fundamental barrier to efficient SLM deployment.

To address this, prior work has explored static vocabulary pruning strategies (Ushio et al., 2023; Yang et al., 2022), which reduce vocabulary size by eliminating rare or irrelevant tokens based on curated corpora. While these approaches are well-motivated, they suffer from key limitations due to their static and coupled design, which assumes a single, globally pruned vocabulary is applied on all vocabulary-related components (i.e., tokenizer, embedding, and LM head). This design introduces two major issues. First, **premature information loss** arises because pruning the tokenizer, embedding, and LM head altogether alters the input representation passed to the transformer. These modified inputs may differ significantly from those seen during pretraining, causing distributional shifts and information loss from the prefill stage. Notably, each pruning step in the pipeline introduces information loss that accumulates, leading to cumulative performance degradation during inference. Second, **lack of flexibility and adaptability** as the static strategy limits the model’s adaptability across diverse tasks. Supporting different task configurations typically requires duplicating multiple copies of vocabulary-related components, resulting in substantial storage overhead and increased complexity in deployment.

Based on empirical observations and theoretical analysis, we derive two key principles for efficient design. **Lexical locality** captures the empirical observations that, in common downstream tasks, generation relies on a highly localized vocabulary, where each output depends on a small subset of input tokens and a limited set of task-specific tokens. **Computation asymmetry** reflects the distinct computational characteristics of the embedding and LM head: the embedding layer mainly leverages lookup operations, which are computationally cheap but memory-bandwidth bound, whereas the LM head is compute-intensive with massive matrix multiplication requiring immense floating-point power and thus better suited for GPUs. Existing pruning methods ignore such asymmetry, missing opportunities for system-level optimization.

Guided by these principles, we propose **VocabTailor**, a flexible and efficient framework for dynamic vocabulary selection. VocabTailor is based on two main pillars. First, we adopt a **decoupled design** for vocabulary-related components. We retain the full tokenizer and offload the embedding layer to CPU memory. Since the embedding lookup is a memory-intensive operation with $\mathcal{O}(1)$ com-

putational complexity, we can strategically offload it to free up valuable GPU memory, with minimal overhead to the overall system performance. Second, we propose **hybrid static-dynamic vocabulary selection**: at runtime, we dynamically select and load input-relevant tokens while maintaining a small, static set of task-specific tokens to ensure stable and efficient computation in the LM head.

This design enables substantial memory savings without compromising input fidelity and model generality. Compared to static pruning, which retains the union of all input and output tokens ($\bigcup \mathcal{I}_i \cup \bigcup \mathcal{O}_i$), VocabTailor only needs $\mathcal{I}_i \cup \mathcal{T}$ at inference time, where \mathcal{I}_i , \mathcal{O}_i are input and output tokens for example i , and \mathcal{T} is a small, fixed task-specific token set. Since $|\bigcup \mathcal{I}_i| \gg |\mathcal{I}_i|$, this leads to substantial memory savings and improved task adaptability.

In summary, our contributions are:

1. We present the first systematic analysis of vocabulary management in LLMs through the lens of *lexical locality* and *computation asymmetry*.
2. We propose **VocabTailor**, a flexible, memory-efficient, and task-adaptive framework that supports a hybrid static-dynamic vocabulary selection strategy along with an enhanced profiling strategy.
3. Across five representative downstream tasks—machine translation, summarization, code completion, information extraction, and math problem solving—VocabTailor reduces the memory usage of vocabulary-related components (i.e., the embedding matrix and LM head) by up to 99%, with minimal or no performance degradation.

2 Related Work

2.1 Small Language Model

Small language models (SLMs) are compact alternatives to large language models (LLMs), which are designed for efficiency, lower computational costs, and deployment on resource-constrained devices. While LLMs like GPT-4 (Achiam et al., 2023), LLaMA (Touvron et al., 2023a,b), Claude (Anthropic, 2024, 2025), Gemini (Team et al., 2024), Qwen (Bai et al., 2023), and DeepSeek (Guo et al., 2025) have achieved widespread success across real-life applications, SLMs are gaining

attention due to their suitability in GPU memory-constrained environments, personal devices, and task-specific scenarios. They offer a scalable, efficient, and sustainable solution tailored for real-time and on-device applications (Lamaakal et al., 2025). With careful selection, small open models can rival and even outperform LLMs while offering improved speed and memory efficiency (Sinha et al., 2024).

Despite advances in architectures, training techniques, and model compression techniques, SLMs still face challenges, including trade-offs between model size and accuracy, generalization limitations, and concerns over bias and privacy (Van Nguyen et al., 2024; Lamaakal et al., 2025). A common approach to building SLMs is distilling them from LLMs while retaining the same tokenizer and vocabulary. This typically causes vocabulary-related components (i.e., embedding and LM Head) to account for a large proportion of the model’s total parameters. This makes vocabulary pruning an effective optimization strategy for SLMs.

2.2 Tokenization

Tokenization is a fundamental preprocessing step in Natural Language Processing (NLP) that splits text into smaller units called tokens (e.g., words or characters), which form the input to downstream tasks. Over the years, various tokenization techniques have been developed (Sennrich et al., 2016; Kudo and Richardson, 2018; Devlin et al., 2019). Among these, Byte-Pair Encoding (BPE) (Sennrich et al., 2016) has become one of the most widely used. Originally developed for data compression, BPE was adapted to tokenize text by iteratively merging the most frequent adjacent symbol pairs starting from individual characters until a target vocabulary size is reached. This approach enables efficient representation of frequent words with fewer tokens while breaking rare or unseen words into informative subword units. BPE’s widespread adoption across transformer architectures has established it as a core component of LLM infrastructure.

State-of-the-art LLMs such as GPT-4 (Achiam et al., 2023), LLaMA (Touvron et al., 2023a,b), Gemini (Team et al., 2024), Claude (Anthropic, 2024, 2025), and DeepSeek (Guo et al., 2025) rely on BPE-based tokenizers to balance vocabulary efficiency and expressiveness. However, the resulting vocabulary sizes are often large, leading to large embedding matrices and LM heads, which increase computational and memory overhead during

inference. This scalability bottleneck has motivated research into vocabulary reduction and adaptive tokenization strategies, especially for resource-constrained deployments.

2.3 Vocabulary Pruning

Vocabulary pruning has emerged as a key area of research in NLP, particularly for scaling and deploying efficient language models. During the BERT era, interest in pruning surged as researchers explored ways to streamline BERT and other transformer-based models. More recently, with the rise of small language models (SLMs), efficient vocabulary selection has once again become a pressing concern.

Vocabulary-trimming (Ushio et al., 2023) and TextPruner (Yang et al., 2022) were proposed to improve efficiency and reduce model size by removing tokens irrelevant to the target language or rarely seen in downstream tasks. Both methods follow a static pruning strategy: they identify language-specific or task-relevant tokens from a curated corpus and prune the vocabulary accordingly. While this process simultaneously reduces the size of the tokenizer, embedding, and LM head, it introduces premature information loss, leading to cumulative performance degradation, and limited flexibility and adaptability, resulting in substantial memory overhead and increased deployment complexity. These two major limitations underscore the need for more flexible and dynamic vocabulary pruning strategies that are task-specific, adaptable, and generalizable across deployment settings. In response, we propose our framework in Section 3.

3 Method

3.1 VocabTailor Framework

VocabTailor (shown in Figure 1) decouples vocabulary management across vocabulary-related components. Unlike static vocabulary pruning that uniformly reduces the tokenizer, embedding layer, and LM head, VocabTailor treats each component based on its unique computational and storage properties.

For the **tokenizer**, we retain the full vocabulary to preserve input fidelity and prevent information loss during the prefill stage. This avoids the cumulative performance degradation caused by pruned tokenizers, which often fail to fully capture the input expressiveness.

For the embedding layer and LM head, despite their similar roles (tokenize and detokenize), a fun-

damental computation asymmetry exists between them. The **embedding layer**, which relies on simple lookup operations, is naturally CPU-friendly and incurs negligible runtime overhead when offloaded to CPU memory from high-cost accelerator memory (Yu et al., 2025). Thus, we retain it fully and offload it to CPU. In resource-constrained deployment settings such as edge devices or platforms with unified CPU–GPU memory, even CPU-resident embeddings may impose non-trivial memory overhead. To address this, VocabTailor supports disk-backed embedding offloading using the Lightning Memory-Mapped Database (LMDB). Embedding vectors are stored as key–value entries and retrieved on demand during inference, preserving full-vocabulary access while further reducing memory footprint. Implementation details and analysis are in Appendix B.

The **LM head**, which performs compute-intensive matrix multiplications, must remain on GPU for efficient inference. To optimize its memory footprint, we introduce a hybrid static-dynamic vocabulary selection strategy. This approach leverages the lexical locality and significantly reduces memory usage while preserving downstream task performance and enabling flexible, task-specific adaptation without model duplication.

This decoupled architecture addresses the key limitations of existing approaches while providing a theoretically grounded and practically efficient solution for large-scale vocabulary management. The detailed design and implementation of vocabulary selection strategy are presented in Sections 3.2–3.4.

3.2 Motivation: Analysis of Lexical Locality

To investigate the essence of efficient vocabulary selection, we analyze input-output pairs across diverse datasets aligned with downstream tasks. Our empirical analysis reveals two fundamental properties of lexical locality that motivate VocabTailor’s dynamic vocabulary selection strategy.

Observation 1 *Input-Driven Locality: Common downstream tasks exhibit strong input-output lexical overlap—each output contains a small subset of input tokens.*

In various downstream tasks, the model output reuses tokens from the input (shown in Figure 2). This phenomenon is particularly pronounced in text extraction tasks (e.g., span-based QA or named entity recognition), where output tokens are typically a subset of the input tokens. In code com-

pletion, generated code often replicates variable names, function names, and other identifiers from the input context. Similarly, in text summarization, the generated summary contains large spans from the source document. For instance, in summarization, on average 61.9% of tokens in generated summaries are copied from the input document. From an information-theoretic perspective, the input context dramatically reduces the entropy of the output token distribution, constraining generation to a much smaller effective vocabulary. Thus, preserving input vocabulary is critical for maintaining performance during vocabulary reduction.

Observation 2 *Task-Driven Locality: Vocabulary required for generation is highly localized—each output depends on a limited set of task-specific tokens.*

Due to input diversity across datasets, the union of all input tokens is significantly larger than the tokens required for any single generation instance. Each input introduces unique tokens, making the aggregate input vocabulary much larger than individual requirements. The remaining output tokens—those not found in the corresponding input—form a relatively small, task-specific set \mathcal{T} that captures the essential generation patterns for the task.

Formally, let \mathcal{I} denote the set of all input tokens, \mathcal{I}_i the set of input tokens for instance i in a dataset, \mathcal{O} the set of all output tokens, and $\mathcal{T} \subset \mathcal{O}$ the set of all task-specific tokens in outputs. We observe:

$$|\mathcal{I}| = \left| \bigcup_i \mathcal{I}_i \right| \gg |\mathcal{I}_i| \text{ and } |\mathcal{O}| > |\mathcal{T}|$$

These observations expose the fundamental inefficiency of static vocabulary pruning that operates on the union $\mathcal{I} \cup \mathcal{O}$ for every example. VocabTailor exploits this lexical locality by dynamically constructing active per-example vocabularies using only per-example input tokens \mathcal{I}_i and a compact task-specific set \mathcal{T} :

$$|\mathcal{I} \cup \mathcal{O}| \gg |\mathcal{I}_i \cup \mathcal{T}|$$

This dynamic targeting of the much smaller $\mathcal{I}_i \cup \mathcal{T}$ enables substantial memory savings without compromising generation quality, as it only retains the tokens necessary for each inference instance.

3.3 The Hybrid Vocabulary Selection Strategy

VocabTailor utilizes lexical locality through a hybrid architecture that combines dynamic runtime

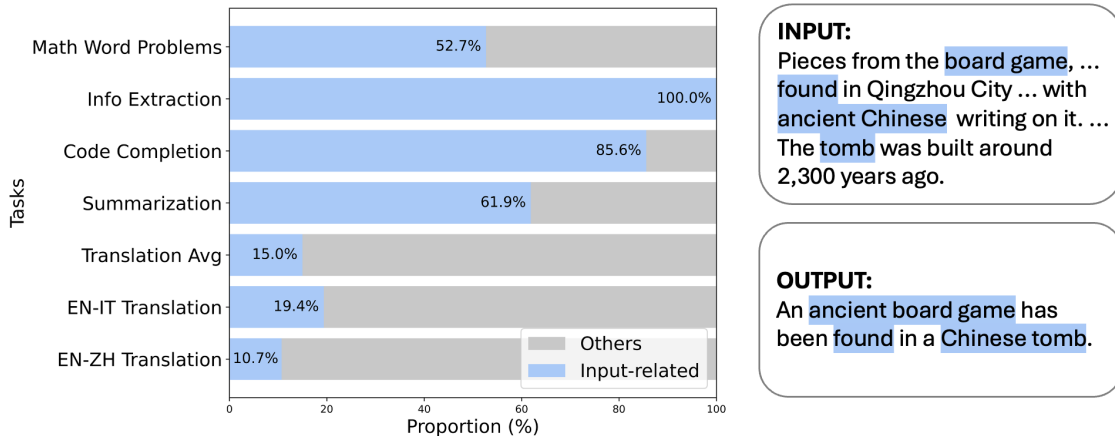


Figure 2: Left: Input-output lexical overlap ratio. Right: Example of lexical overlap in a summarization task.

adaptation with static offline optimization, enabling efficient vocabulary management without sacrificing generation flexibility.

3.3.1 Dynamic Selection (Runtime Behavior)

At the start of each inference instance, VocabTailor identifies the unique input token set \mathcal{I}_i from the input text and selectively loads the corresponding LM head weight vectors from CPU to GPU. This selective loading exploits the input-driven locality principle, ensuring that only input-relevant vocabulary components are active during generation.

Beyond selecting input-related tokens at runtime, VocabTailor must efficiently extend the LM head on the GPU. A naive implementation that repeatedly concatenates static and dynamic weights and instantiates new linear layers incurs unnecessary CPU-GPU transfers and module creation overhead. To mitigate this, VocabTailor supports PreAlloc, a pre-allocation strategy that leverages the structure of input-related tokens to amortize LM head construction across inputs.

To minimize the frequency of linear module creation, we investigate the number of unique input-related tokens across different tasks, and we observe that it is relatively small and the distribution is concentrated. Based on this insight, we pre-allocate a small memory space for the dynamic LM head weight. As long as the dynamic weights fit within the buffer size, no concatenation and no new linear module creation are required. The static weights remain on the GPU without movement, so the memory movement overhead is minimized. In the infrequent event that the size of dynamic weights exceeds the current buffer size, a new Linear module with a larger buffer size is created in a

way similar to the naive implementation, thereby significantly reducing the expected latency overhead. Implementations and performance evaluation are detailed in Appendix A.

3.3.2 Static Selection (Offline Construction)

The static component maintains a compact, task-specific core vocabulary \mathcal{T} that captures essential output tokens independent of input context. While simple frequency-based filtering proves inadequate due to noise from typos and multilingual interference in broad corpora, we introduce a theoretically grounded filtering pipeline (Algorithm 1) that constructs \mathcal{T} through fine-grained analysis.

3.4 Fine-grained Construction of the Static Task Vocabulary

Our static vocabulary construction algorithm (Algorithm 1) addresses the challenge of isolating task-essential tokens via a three-stage filtering process: **Input-Aware Filtering**. We first exclude all input tokens from the candidate vocabulary, isolating tokens the model must generate without input cues (e.g., function keywords in code generation, discourse markers in summarization). This step directly implements task-driven locality by identifying the irreducible core \mathcal{T} that cannot be derived from the input context.

Language-Specific Filtering. To suppress noise in mixed-language scenarios (e.g., code datasets with multilingual comments), we apply Unicode block analysis to retain only tokens belonging to the target language family. This heuristic-based approach effectively handles cross-lingual interference while preserving task-relevant vocabulary.

Tolerance Filtering. When additional vocabu-

lary reduction is required, we formulate vocabulary selection similar to a pruning problem. We define the impact metric $df(v)$ as the document frequency—the fraction of instances where token v appears in the ground truth. Tokens are sorted by ascending $df(v)$ and iteratively removed until the cumulative document frequency of pruned tokens reaches a user-defined tolerance threshold τ . This threshold bounds the worst-case performance drop, where $\tau = 0.01$ ensures that at most 1% of profiling samples lose a critical token.

VocabTailor’s hybrid approach bridges theoretical rigor (exploiting lexical locality principles) with practical efficiency (optimizing GPU-CPU memory hierarchy). The static-dynamic decomposition enables substantial memory savings while preserving the model’s generational capabilities, providing a scalable solution for deploying SLMs across diverse applications with controllable efficiency-accuracy trade-offs.

4 Experiments

4.1 Settings

4.1.1 Tasks and datasets

We evaluate VocabTailor on five representative SLM downstream tasks across diverse domains: machine translation, summarization, code completion, information extraction, and math word problem solving. For machine translation, we involve English-to-Italian and English-to-Chinese translation, as Chinese is logographic with minimal morphology and a large character set, while Italian is alphabetic and morphologically rich. We use the Opus-100 corpus (Zhang et al., 2020) as a profiling dataset and WMT24++ (Deutsch et al., 2025) for the evaluation. For summarization, we use the XSum training set (Narayan et al., 2018) to profile and evaluate on its test set. For code completion, the CodeContests+ corpus (Wang et al., 2025) is used for profiling, while evaluation is conducted on SAFIM (Gong et al., 2024). For information extraction, we use the SQuAD (Rajpurkar et al., 2016) training set for profiling and its test set for evaluation. In math problem solving, GSM8K (Cobbe et al., 2021) served as the profiling corpora, with evaluation performed on MAWPS (Koncel-Kedziorski et al., 2016). Each dataset is selected for its strong alignment with the target task.

4.1.2 Evaluation metrics

We include sacreBLEU (Post, 2018), METEOR (Banerjee and Lavie, 2005), and COMET (Rei et al., 2020) for machine translation. For summarization, we use Rouge-1, Rouge-2, and Rouge-L scores (Lin, 2004). Pass@1 (Chen et al., 2021) is used for code completion. We use F1 score for information extraction, and accuracy for math problem solving. These metrics are standard in the field and provide robust measures of model performance across the target tasks. For efficiency evaluation, we report Time to First Token (TTFT), Time Per Output Token (TPOT), Token Per Second (TPS), and Peak VRAM usage.

4.1.3 Models

For machine translation, we use Qwen3-1.7B (Yang et al., 2025). We employ Llama 3.2 3B for summarization and Llama 3.2 1B for information extraction (Dubey et al., 2024). For summarization, the base model is fine-tuned for a better base performance. For code completion, we choose deepseek-coder-1.3b-base (Guo et al., 2024), and for math problem solving, we apply rho-math-1b-interpreter-v0.1 (Lin et al., 2024). For efficiency evaluation, we use Qwen3-0.6B (Yang et al., 2025).

4.1.4 Baselines and Other Settings

We compare our method (VocabTailor) with the original model (Original) and static vocabulary pruning (VP). Static vocabulary pruning follows the common routines of corpus-based filtering. For a fair comparison, VP and VocabTailor use the same profiling corpora. We set the tolerance threshold $\tau = 0.01$ on all tasks for VocabTailor. As a hybrid dynamic-static framework, the vocabulary size of VocabTailor model varies in each single inference. Here we report the average vocabulary size for each downstream task. For all models, we set the temperature to 0 to avoid randomness. For efficiency evaluation, we use the Qwen3-0.6B model along with 100 prompts from the machine translation (English-to-Chinese) task. We report the latency and VRAM usage on multiple devices, including NVIDIA A100 GPU, Apple Silicon M1 Pro, and Jetson Orin Nano Super.

4.2 Results

As shown in Table 1, VocabTailor consistently achieves substantial vocabulary reduction while maintaining or even improving task performance

Task	Model	Vocabulary	(% Original)	Metric
Machine Translation English→Chinese	Original	151,643		13.48/12.00/81.16
	VP	64,117	(42.28%)	14.84/11.84/81.19
	VocabTailor (Ours)	18,874 + [40]	(12.47%)	15.39/12.69/81.44
Machine Translation English→Italian	Original	151,643		24.33/48.95/73.87
	VP	65,518	(43.21%)	22.59/46.74/74.93
	VocabTailor (Ours)	24,185 + [14]	(15.96%)	21.13/46.97/75.49
Summarization	Original	128,000		0.36/0.15/0.29
	VP	59,613	(46.57%)	0.36/0.15/0.29
	VocabTailor (Ours)	36,332 + [15]	(28.40%)	0.36/0.15/0.29
Code Completion	Original	32,000		54.10%
	VP	24,888	(77.78%)	8.12%
	VocabTailor (Ours)	3,521 + [58]	(11.18%)	53.87%
Information Extraction	Original	128,000		38.40
	VP	49,106	(38.36%)	2.58
	VocabTailor (Ours)	[105]	(0.08%)	62.73
Math Problem Solving	Original	32,000		88.40
	VP	10,300	(32.19%)	87.80
	VocabTailor (Ours)	5,135 + [14]	(16.09%)	88.40

Table 1: Results on machine translation (sacreBLEU/METEOR/COMET), summarization (Rouge-1/Rouge-2/Rouge-L), code completion (Pass@1), information extraction (F-1), and math problem solving (Accuracy), including the vocabulary size and the ratio to the original model (%). For VocabTailor, vocabulary size consists of task-specific tokens and the average number of dynamic tokens highlighted in brackets. The best results are in bold characters.

compared with the original model, outperforming static pruning in nearly all cases.

4.2.1 Qualitative Evaluation

Machine Translation. In English-to-Chinese translation, VocabTailor achieves the best results (SacreBLEU 15.39, METEOR 12.69, and COMET 81.44), while using only 12% of the full vocabulary. These improvements are notable given the inherent difficulty of vocabulary pruning in high-character-set languages like Chinese. The results suggest that VocabTailor’s dynamic vocabulary selection, which retains input-relevant tokens while filtering low-utility output tokens, preserves both surface-level fluency and semantic adequacy. Trimming irrelevant vocabulary reduces noise in the output distribution and improves decoding stability.

In English-to-Italian translation, VocabTailor attains the highest COMET (75.49), indicating strong preservation of semantic and contextual meaning. We attribute this to its ability to retain morphology-relevant tokens while removing unrelated vocabulary, which is particularly beneficial for morphologically rich languages such as Italian. It retains only

16% of the original vocabulary, yet performs competitively on SacreBLEU (21.13) and METEOR (46.97), with modest drops compared to the unpruned model. In contrast, VP removes less vocabulary (43%) but still causes moderate degradation across all metrics, suggesting that static pruning may eliminate rare but task-relevant tokens. VocabTailor’s ability to adapt to Italian’s rich inflectional morphology further reinforces the general applicability of the method.

Summarization. All three models yield identical ROUGE scores, indicating no difference in summary quality. VocabTailor achieves this with just 28% of the original vocabulary, 18% fewer tokens than VP. These results demonstrate that aggressive vocabulary reduction is possible without compromising output quality when applying a hybrid dynamic-static vocabulary selection.

Code Completion. On the SAFIM benchmark with deepseek-coder-1.3b-base, VP causes a dramatic performance drop: Pass@1 falls from 54.10% to 8.12%, despite a modest vocabulary reduction to 78% of the original size. This highlights a key limitation of corpus-based pruning for code:

essential elements such as variable names and identifiers may be discarded if they appear infrequently in the profiling corpus. In contrast, VocabTailor retains only 11% of the original size and achieves a high Pass@1 of 53.87%. These findings underscore the robustness of our input- and task-specific vocabulary selection strategy for generation-intensive tasks such as code synthesis.

Information Extraction. For information extraction tasks, we evaluate on the SQuAD dataset using Llama 3.2 1B. VocabTailor achieves the most striking result: with just 0.08% of the original vocabulary retained (and no static tokens at all), it attains an F1 score of 62.73, outperforming both the Original (38.40) and VP (2.58) by large margins. This improvement is expected for extractive tasks, where the output vocabulary is typically a subset of the input. By dynamically preserving all input tokens, VocabTailor guarantees coverage of the necessary output vocabulary, whereas VP may prune such tokens during static filtering. This highlights the limitation of static, corpus-based pruning under high input–output overlap and the advantage of VocabTailor in this setting.

Math Word Problem Solving. The unpruned model achieves a score of 88.40, which is fully preserved under VocabTailor, even though it reduces the vocabulary to just 16% of the Original. VP, while also reducing vocabulary size to 32%, causes a slight performance drop to 87.80. This shows that even in tasks requiring high-precision symbolic handling, our approach remains effective.

Across five tasks, VocabTailor consistently demonstrates strong performance while substantially reducing vocabulary size. In many cases, it matches or even exceeds the performance of the unpruned model. Compared to static pruning, which often compromises accuracy, VocabTailor reduces vocabulary more aggressively (up to 99%, e.g., to as low as 0.08% of the Original in information extraction) without sacrificing model quality. These findings support that dynamic vocabulary selection is a practical and efficient approach to deploying SLMs in resource-constrained environments.

4.2.2 Efficiency Evaluation

We evaluate the inference efficiency of VocabTailor compared to the original small language model across diverse hardware platforms, including edge devices (Apple Silicon M1 Pro and Jetson Orin Nano Super) and a high-end GPU (NVIDIA A100). Table 2 shows that VocabTailor consistently re-

duces peak VRAM usage by 22-23% (from 1.11 GB to 0.85 GB on M1 Pro chip and 1.17 GB to 0.91 GB on NVIDIA GPUs), highlighting its effectiveness in alleviating memory bottlenecks on resource-constrained devices. In the decoding phase, VocabTailor achieves modest improvements in Time Per Output Token (TPOT) on lower-power hardware (up to 17% faster on M1 Pro) with negligible impact on the A100, owing to reduced computation in the dynamically trimmed LM head. Time to First Token (TTFT) increases on all devices, due to overhead from on-demand vocabulary loading during the prefill phase. However, as the runtime is dominated by decoding, VocabTailor has comparable TPS to the original model.

4.3 Ablation Study

To understand the contributions of each component in our framework, we conduct a series of ablation experiments on the SAFIM dataset using the deepseek-coder-1.3b-base model for the code completion task. We primarily focus on evaluating the impact of different vocabulary configurations, including the static and dynamic components, as well as our proposed three-stage filtering process.

Impact of Dynamic vs. Static Vocabulary Components. We first evaluate the impact of the dynamic and static vocabulary components, both individually and in combination. As shown in Table 3, using only the dynamic part that contains tokens profiled from the specific input examples results in a significant performance drop (Pass@1 of 36.06%). This demonstrates that input tokens alone lack the broader coverage needed for robust code completion. Using only the static part (task-specific tokens) achieves a Pass@1 of 52.30%. However, this still underperforms the full static–dynamic configuration, which achieves the best result at 53.87% with nearly the same vocabulary size. This indicates that while the static tokens carry most of the task-relevant capacity, including the dynamic tokens adds critical input-specific nuances, and their combination is essential for optimal performance.

We also compare our static-only approach with VP: VP retains 78% of the original vocabulary—more than our approach—it results in a drastically lower Pass@1 of 8.12%. This stark contrast underscores a key insight: while both VP and our static-only setup are static, VP’s direct modification of the tokenizer and embeddings damages input representations, leading to severe degradation. In contrast, our method retains the full tokenizer and embed-

Model	Device	Avg. TTFT (ms)	Avg. TPOT (ms)	TPS (token/s)	Peak VRAM (GB)
Original	Apple Silicon M1 Pro	23.58	75.10	18.82	1.11
	Jetson Orin Nano Super	4.40	130.80	7.68	1.17
	NVIDIA A100	1.74	24.21	41.86	1.17
VocabTailor	Apple Silicon M1 Pro	25.73	62.03	20.27	0.85
	Jetson Orin Nano Super	20.62	127.22	7.93	0.91
	NVIDIA A100	3.71	24.30	41.65	0.91

Table 2: Performance comparison between the Original model and VocabTailor across different devices.

Model	Vocabulary Pass@1
Original	100% 54.10%
VP	77.80% 8.12%
Dynamic + Static ($\tau = 0.01$)	11.18% 53.87%
Dynamic only	0.81% 36.06%
Static only ($\tau = 0.01$)	11.00% 52.30%

Table 3: Comparison of the dynamic and static components in VocabTailor.

ding, thereby preserving representational integrity and maintaining high performance even with significantly fewer tokens.

Impact of Input-aware and Language-specific Filtering. As discussed earlier, VocabTailor calibrates the static vocabulary through three-stage filtering. In Table 4, starting from the unfiltered static vocabulary (78% of the Original), applying input-aware filtering (IA) reduces the size to 53% with virtually no performance loss (54.07% vs. 54.06%). Adding language-specific filtering (LS) further reduces the vocabulary size to 46%, while performance slightly improves to 54.09%. This improvement likely stems from the removal of noisy or irrelevant tokens from multilingual corpora, allowing the model to focus on task-relevant representations. These results demonstrate that IA and LS can significantly compress the vocabulary without degrading accuracy, validating the effectiveness of our static token selection process.

Impact of Tolerance Filtering. Lastly, we analyze the effect of tolerance filtering, which allows for further reduction of rarely activated tokens based on cumulative document frequency. In Table 5, we vary the tolerance threshold (τ) to observe the trade-off between vocabulary size and accuracy. At $\tau = 0$, we preserve all profiled tokens, achieving a Pass@1 of 54.09%. As tolerance increases, more

Model	Vocabulary Pass@1
Original	100.00% 54.10%
Dynamic + Unfiltered Static	77.78% 54.07%
Dynamic + IA	52.85% 54.06%
Dynamic + IA + LS	45.61% 54.09%

Table 4: Comparison of input-aware filtering (IA) and language-specific filtering (LS) in VocabTailor.

tokens are filtered out, reducing the vocabulary to as low as 2.5% of the original size ($\tau = 0.10$), with a gradual decline in performance. Importantly, even with $\tau = 0.01$, the vocabulary shrinks to 11% with only 1.8% drop in Pass@1, suggesting that our method is robust to aggressive reduction. This highlights the flexibility of tolerance as a tuning knob to balance compression vs. accuracy.

Model	Vocabulary	Pass@1
$\tau = 0$	45.61%	54.09%
$\tau = 0.01$	11.18%	53.87%
$\tau = 0.02$	7.28%	53.71%
$\tau = 0.10$	2.50%	52.28%

Table 5: Comparison of different tolerance thresholds in VocabTailor.

5 Conclusion

In this paper, we propose a flexible and efficient vocabulary selection framework to reduce memory usage during SLM inference. By identifying and leveraging lexical locality and computation asymmetry, our method reduces up to 99% of the memory associated with vocabulary-related components (i.e., the embedding matrix and LM head) of SLM while maintaining performance on representative downstream tasks.

Limitations

The proposed framework presented in this paper is only explored on language downstream tasks of SLMs. This method may be extended and applied to VLMs, ALMs, and MLLMs for memory-efficient inference on image/video understanding and audio generation tasks. The method focuses on the memory reduction of SLMs, so (1) although our method mildly increases model decoding speed, inference acceleration is not a goal of this work, and (2) while the method can be applied to LLM, the reduction in memory will be limited.

Acknowledgments

This work is supported by the JC STEM Lab of AI for Science and Engineering, funded by The Hong Kong Jockey Club Charities Trust and the Research Grants Council of Hong Kong (Project No. CUHK14213224).

References

- Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, and 1 others. 2023. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774*.
- Anthropic. 2024. The Claude 3 Model Family: Opus, Sonnet, Haiku. Technical report.
- Anthropic. 2025. Introducing Claude 4. <https://www.anthropic.com/news/claude-4>.
- Jinze Bai, Shuai Bai, Yunfei Chu, Zeyu Cui, Kai Dang, Xiaodong Deng, Yang Fan, Wenbin Ge, Yu Han, Fei Huang, and 1 others. 2023. Qwen technical report. *arXiv preprint arXiv:2309.16609*.
- Satanjeev Banerjee and Alon Lavie. 2005. **METEOR: An automatic metric for MT evaluation with improved correlation with human judgments**. In *Proceedings of the ACL Workshop on Intrinsic and Extrinsic Evaluation Measures for Machine Translation and/or Summarization*, pages 65–72, Ann Arbor, Michigan. Association for Computational Linguistics.
- Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, and 1 others. 2020. Language models are few-shot learners. *Advances in neural information processing systems*, 33:1877–1901.
- Mark Chen, Jerry Tworek, Heewoo Jun, Qiming Yuan, Henrique Ponde De Oliveira Pinto, Jared Kaplan, Harri Edwards, Yuri Burda, Nicholas Joseph, Greg Brockman, and 1 others. 2021. Evaluating large language models trained on code. *arXiv preprint arXiv:2107.03374*.
- Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, Mark Chen, Heewoo Jun, Lukasz Kaiser, Matthias Plappert, Jerry Tworek, Jacob Hilton, Reiichiro Nakano, and 1 others. 2021. Training verifiers to solve math word problems. *arXiv preprint arXiv:2110.14168*.
- Daniel Deutsch, Eleftheria Briakou, Isaac Caswell, Mara Finkelstein, Rebecca Galor, Juraj Juraska, Geza Kovacs, Alison Lui, Ricardo Rei, Jason Riesa, and 1 others. 2025. Wmt24++: Expanding the language coverage of wmt24 to 55 languages & dialects. *arXiv preprint arXiv:2502.12404*.
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. **BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding**. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 4171–4186, Minneapolis, Minnesota. Association for Computational Linguistics.
- Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, and 1 others. 2024. The llama 3 herd of models. *arXiv e-prints*, pages arXiv–2407.
- Linyuan Gong, Sida Wang, Mostafa Elhoushi, and Alvin Cheung. 2024. Evaluation of llms on syntax-aware code fill-in-the-middle tasks. In *International Conference on Machine Learning*, pages 15907–15928. PMLR.
- Daya Guo, Dejian Yang, Haowei Zhang, Junxiao Song, Ruoyu Zhang, Runxin Xu, Qihao Zhu, Shitong Ma, Peiyi Wang, Xiao Bi, and 1 others. 2025. Deepseek-r1: Incentivizing reasoning capability in llms via reinforcement learning. *arXiv preprint arXiv:2501.12948*.
- Daya Guo, Qihao Zhu, Dejian Yang, Zhenda Xie, Kai Dong, Wentao Zhang, Guanting Chen, Xiao Bi, YK Li, and 1 others. 2024. Deepseek-coder: When the large language model meets programming—the rise of code intelligence. *arXiv preprint arXiv:2401.14196*.
- Rik Koncel-Kedziorski, Subhro Roy, Aida Amini, Nate Kushman, and Hannaneh Hajishirzi. 2016. Mawps: A math word problem repository. In *Proceedings of the 2016 conference of the north american chapter of the association for computational linguistics: human language technologies*, pages 1152–1157.
- Taku Kudo and John Richardson. 2018. Sentencepiece: A simple and language independent subword tokenizer and detokenizer for neural text processing. *arXiv preprint arXiv:1808.06226*.

- Ismail Lamaakal, Yassine Maleh, Khalid El Makkaoui, Ibrahim Ouahbi, Paweł Pławiak, Osama Alfarraj, May Almousa, and Ahmed A Abd El-Latif. 2025. Tiny language models for automation and control: Overview, potential applications, and future research directions. *Sensors*, 25(5):1318.
- Chin-Yew Lin. 2004. Rouge: A package for automatic evaluation of summaries. In *Text summarization branches out*, pages 74–81.
- Zhenghao Lin, Zhibin Gou, Yeyun Gong, Xiao Liu, Yelong Shen, Ruochen Xu, Chen Lin, Yujiu Yang, Jian Jiao, Nan Duan, and 1 others. 2024. Rho-1: Not all tokens are what you need. *arXiv preprint arXiv:2404.07965*.
- Shashi Narayan, Shay B. Cohen, and Mirella Lapata. 2018. Don't Give Me the Details, Just the Summary! Topic-Aware Convolutional Neural Networks for Extreme Summarization. *Preprint*, arXiv:1808.08745.
- Matt Post. 2018. A call for clarity in reporting BLEU scores. In *Proceedings of the Third Conference on Machine Translation: Research Papers*, pages 186–191, Belgium, Brussels. Association for Computational Linguistics.
- Pranav Rajpurkar, Jian Zhang, Konstantin Lopyrev, and Percy Liang. 2016. SQuAD: 100,000+ Questions for Machine Comprehension of Text. *Preprint*, arXiv:1606.05250.
- Ricardo Rei, Craig Stewart, Ana C Farinha, and Alon Lavie. 2020. Comet: A neural framework for mt evaluation. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 2685–2702.
- Rico Sennrich, Barry Haddow, and Alexandra Birch. 2016. Neural Machine Translation of Rare Words with Subword Units. In *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 1715–1725, Berlin, Germany. Association for Computational Linguistics.
- Neelabh Sinha, Viniya Jain, and Aman Chadha. 2024. Are small language models ready to compete with large language models for practical applications? *arXiv preprint arXiv:2406.11402*.
- Gemini Team, Petko Georgiev, Ving Ian Lei, Ryan Burnell, Libin Bai, Anmol Gulati, Garrett Tanzer, Damien Vincent, Zhufeng Pan, Shibo Wang, and 1 others. 2024. Gemini 1.5: Unlocking multimodal understanding across millions of tokens of context. *arXiv preprint arXiv:2403.05530*.
- Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, and 1 others. 2023a. Llama: Open and efficient foundation language models. *arXiv preprint arXiv:2302.13971*.
- Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shrutu Bhosale, and 1 others. 2023b. Llama 2: Open foundation and fine-tuned chat models. *arXiv preprint arXiv:2307.09288*.
- Asahi Ushio, Yi Zhou, and Jose Camacho-Collados. 2023. Efficient Multilingual Language Model Compression through Vocabulary Trimming. In *Findings of the Association for Computational Linguistics: EMNLP 2023*, pages 14725–14739, Singapore. Association for Computational Linguistics.
- Chien Van Nguyen, Xuan Shen, Ryan Aponte, Yu Xia, Samyadeep Basu, Zhengmian Hu, Jian Chen, Mihir Parmar, Sasidhar Kunapuli, Joe Barrow, and 1 others. 2024. A survey of small language models. *arXiv preprint arXiv:2410.20011*.
- Zihan Wang, Siyao Liu, Yang Sun, Hongyan Li, and Kai Shen. 2025. Codecontests+: High-quality test case generation for competitive programming. *arXiv preprint arXiv:2506.05817*.
- An Yang, Anfeng Li, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chang Gao, Chengen Huang, Chenxu Lv, and 1 others. 2025. Qwen3 technical report. *arXiv preprint arXiv:2505.09388*.
- Ziqing Yang, Yiming Cui, and Zhigang Chen. 2022. TextPruner: A Model Pruning Toolkit for Pre-Trained Language Models. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics: System Demonstrations*, pages 35–43, Dublin, Ireland. Association for Computational Linguistics.
- Da Yu, Edith Cohen, Badih Ghazi, Yangsibo Huang, Pritish Kamath, Ravi Kumar, Daogao Liu, and Chiyuan Zhang. 2025. Scaling Embedding Layers in Language Models. *Preprint*, arXiv:2502.01637.
- Biao Zhang, Barry Haddow, and Alexandra Birch. 2023. Prompting Large Language Model for Machine Translation: A Case Study. *Preprint*, arXiv:2301.07069.
- Biao Zhang, Philip Williams, Ivan Titov, and Rico Sennrich. 2020. Improving massively multilingual neural machine translation and zero-shot translation. *arXiv preprint arXiv:2004.11867*.

A Alternative Designs for Dynamic LM Head Construction

The naive (or vanilla) implementation of VocabTailor (Figure 3, left panel) dynamically constructs the LM head by selecting the input-related weights, concatenating them with the task-related static weights, unloading the weights from CPU to GPU, and finally creating a new Linear module for the concatenated LM head. This approach has two limitations. First, to avoid the VRAM peak, each time the concatenation needs to be done on the CPU. This results in the frequent CPU-GPU movement of the task-related static part of the LM head weights, which is highly costly for tasks requiring multi-round interactions or multiple calls to SLMs. Second, we need to initialize a new Linear module each time the dynamic LM head is created, and the creation of a linear module introduces extra latency overhead during the prefill stage. To alleviate these latencies and overheads, we propose two alternative approaches: **SplitLinear** and **PreAlloc**.

A.1 SplitLinear

In the SplitLinear design (Figure 3, middle panel), the static and dynamic parts of the LM head form Linear modules ($M1$ and $M2$) separately. Since the creation of the linear module of the static part can be done in the initialization process, only the creation of the dynamic part needs to be processed during the prefill stage. This design avoids weight concatenation, and the static part of the LM Head can be initialized on the GPU before inference. During the inference process, when the LM head is called, the input is passed through $M1$ and $M2$ independently, and their respective outputs are concatenated to form the final logit. Due to the nature of matrix multiplication, the forward of $M1$ and $M2$ does not interfere with each other, but two GEMMs are required per forward pass.

A.2 PreAlloc

To minimize the frequency of linear module creation, we investigate the number of unique input-related tokens across different downstream tasks. We observe that the number of unique input-related tokens is relatively small, and the distribution is concentrated. Among all five downstream tasks (machine translation, summarization, code completion, information extraction, and math problem solving), the average number of new unique input-related tokens per request remains below 128.

Based on this insight, we pre-allocate a small memory space for the dynamic LM head weight. As shown in the right panel of Figure 3, we created a Linear module on the GPU in advance, where the weight tensor consists of two parts. The first part is the static weights related to the task, which are loaded before inference. The rest of the weight tensor is a zero-initialized buffer. For each input, the dynamic weights related to the input are copied into the buffer. As long as the dynamic weights fit within the buffer size, no concatenation and no new linear module creation are required. In this process, the static weights remain on the GPU without movement, so the memory movement overhead is minimized. In the infrequent event that the size of dynamic weights exceeds the current buffer size, a new Linear module with a larger buffer size is created in a way similar to VocabTailor’s naive implementation, thereby significantly reducing the expected latency overhead.

A.3 Comparison of Dynamic LM Head Construction Strategies

Both strategies are included as inference options in VocabTailor because they offer different advantages and disadvantages tailored to various downstream tasks. The SplitLinear approach is well-suited for downstream tasks with diverse input-related vocabularies but relatively short outputs, such as the information retrieval task, because it provides instant updates and consistent latency regardless of the number of dynamic input-related tokens. However, multiple GEMMs are performed in the forward pass. Conversely, the PreAlloc approach is ideal for downstream tasks with relatively small and stable input-related vocabularies, such as code completion or math problem solving, where the output is primarily composed of task-related terms and only a small portion of dynamic input-related vocabularies like variable names. In such scenarios, pre-allocation effectively eliminates both memory movement and module initialization overhead.

A.3.1 Experimental Setup

We test our proposed methods in the machine translation (English-to-Chinese) task using a subset of 100 examples from the WMT24++ dataset (Deutsch et al., 2025). We use the Qwen3-0.6B (Yang et al., 2025) as the base model. All experiments are conducted on a single NVIDIA A100 GPU. We compared four configurations: (1) Original, the original model with standard full-

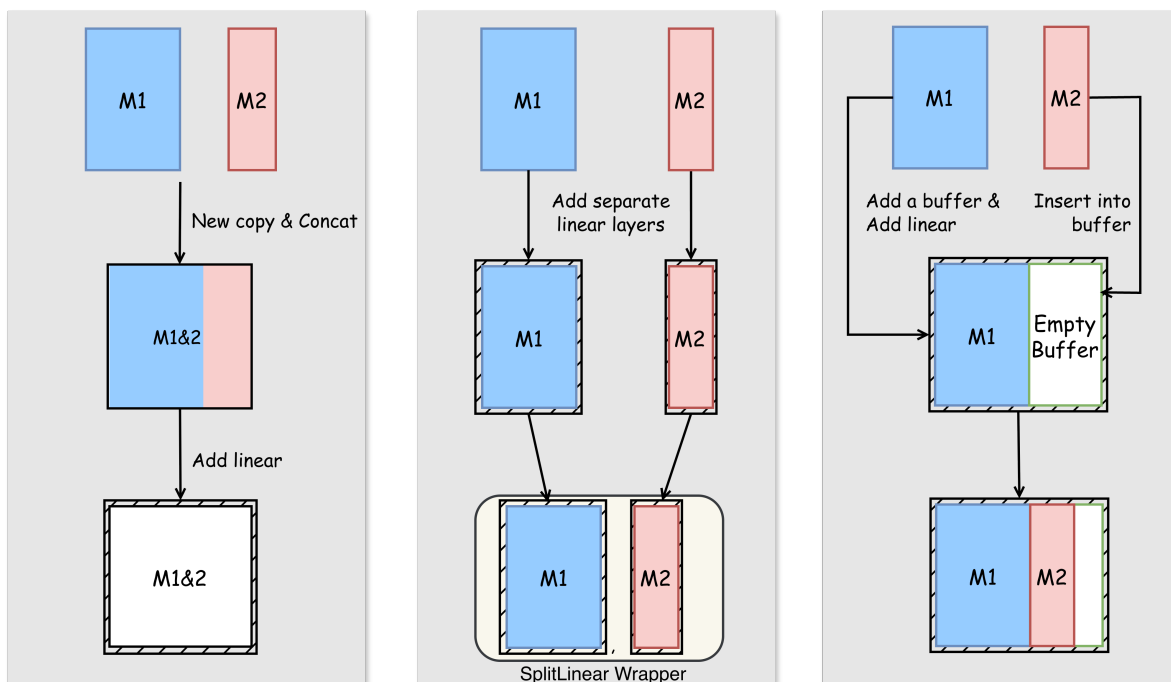


Figure 3: Comparison of dynamic LM head construction workflows. Left panel (Vanilla): The naive re-allocation approach where weights are concatenated and a new Linear module is initialized per request, incurring high movement and initialization costs. Middle panel (SplitLinear): A decoupled architecture where static and dynamic weights form independent Linear modules (M_1 , M_2), allowing the static part to be pre-initialized on the GPU. Right panel (PreAlloc): Our optimized strategy using a pre-warmed GPU buffer. Static weights remain stationary while dynamic weights are copied into a zero-initialized buffer, eliminating module re-initialization and minimizing memory movement overhead.

vocabulary; (2) VT (Vanilla), the VocabTailor naive implementation with dynamic CPU-GPU weight concatenation; (3) VT (SplitLinear), which utilizes dual GEMM operations to bypass concatenation; and (4) VT (PreAlloc), which employs a pre-warmed GPU buffer of size 128. Performance is measured across 100 prompts to capture statistical distributions of latency metrics including Time to First Token (TTFT), Time Per Output Token (TPOT), and end-to-end (E2E) latency. We also report throughput metrics including Token Per Second (TPS), Request Per Second (RPS), and average output length. Memory usage is monitored via the PyTorch CUDA memory management interface to report peak VRAM consumption.

A.3.2 Results

Tables 6 and 7 compare the original model against three VocabTailor (VT) variants, revealing clear trade-offs between initialization overhead, steady-state decoding efficiency, and memory utilization.

Latency behavior. The Original achieves the lowest Time to First Token (TTFT), with a mean of 1.74 ms, reflecting its fully static vocabulary and

absence of runtime adaptation. In contrast, the Vanilla VT incurs a substantial TTFT increase (mean 124.29 ms), which can be attributed to dynamic vocabulary construction and associated runtime bookkeeping. Notably, this overhead is largely confined to the prefill phase: Time Per Output Token (TPOT) remains comparable across all models (≈ 24 ms), indicating that VT does not degrade steady-state decoding once generation begins.

Among different dynamic head construction approaches, PreAlloc is particularly effective at mitigating TTFT overhead, reducing mean TTFT to 3.71 ms, within the same order of magnitude as Original, while preserving the benefits of VT. Tail latency analysis further supports this observation: PreAlloc substantially improves P99 TTFT (17.20 ms vs. 137.43 ms for Vanilla VT), suggesting improved stability under varying prompt complexities. In contrast, SplitLinear does not materially reduce TTFT, as it still requires the dynamic creation of the module M_2 .

End-to-end latency. Despite large TTFT differences, end-to-end (E2E) latency remains broadly

Model	TTFT (ms)				TPOT (ms)				E2E Latency (s)			
	Mean	P50	P90	P99	Mean	P50	P90	P99	Mean	P50	P90	P99
Original	1.74	0.79	0.90	2.07	24.21	23.64	24.09	25.17	1.57	1.56	2.44	4.12
VT (Vanilla)	124.29	123.22	124.63	137.43	23.99	23.47	23.60	24.74	1.66	1.66	2.47	3.91
VT (SplitLinear)	123.86	123.38	124.55	129.62	24.52	24.05	24.19	28.87	1.69	1.69	2.53	3.98
VT (PreAlloc)	3.71	2.63	3.35	17.20	24.30	23.82	24.32	25.96	1.56	1.57	2.38	3.91

Table 6: Latency distribution for different dynamic LM head construction approaches on a machine translation (English-to-Chinese) task ($N = 100$). We report Time to First Token (TTFT), Time Per Output Token (TPOT), and end-to-end (E2E) latency. Percentiles (Mean, P50, P90, P99) are included to characterize tail latency and model stability under varying prompt complexity.

Model	TPS (token/s)	RPS (req/s)	Avg. Output Length	Peak VRAM (GB)
Original	41.86	0.64	65.73	1.17
VT (Vanilla)	42.31	0.56	64.84	0.91
VT (SplitLinear)	41.34	0.55	64.84	0.91
VT (PreAlloc)	41.65	0.55	64.84	0.91

Table 7: Throughput and resource utilization for different dynamic LM head construction approaches on a machine translation (English-to-Chinese) task ($N = 100$). Decode throughput is measured in Tokens Per Second (TPS), computed as total generated tokens divided by total decoding time. Request throughput (RPS), average output length, and peak GPU memory consumption are also reported.

similar across models. All VT variants exhibit mean E2E latency within 1.56–1.69 s, comparable to the Original (1.57 s). This indicates that, for typical MT tasks with non-trivial output lengths, TTFT overhead is amortized over decoding, and overall user-perceived latency is dominated by generation rather than initialization.

Latency decomposition. Figure 4 illustrates the breakdown of total E2E latency consumed by the prefill versus decoding stages. In the Vanilla and SplitLinear implementations, the prefill stage—which is effectively instantaneous in the baseline (0.11%)—surges to over 7.3% of the total execution time. By utilizing a pre-warmed GPU buffer, the pre-allocation strategy successfully reduce the prefilling time, returning the prefill stage to just 0.24% of total latency.

Throughput and resource efficiency. As shown in Table 7, decode throughput (TPS) remains effectively consistent across all models (41–42 tokens/s), confirming that VT and its extensions do not introduce steady-state performance regressions. Request-level throughput, measured by Request Per Second (RPS), is slightly lower for VT-based models, which aligns with their increased per-request initialization cost. In terms of memory footprint, VT variants consistently reduce peak VRAM usage (0.91 GB vs. 1.17 GB for Original),

validating the effectiveness of vocabulary decoupling and offloading.

Overall trade-offs. Taken together, these results highlight that VocabTailor introduces a clear TTFT–memory trade-off: substantial VRAM savings with minimal impact on decode throughput, at the cost of higher initialization latency. Among the different head construction strategies, PreAlloc offers the most favorable balance, largely eliminating TTFT penalties while preserving VRAM savings and decode efficiency. SplitLinear while maintaining comparable throughput and memory characteristics, provides limited benefit for latency-sensitive scenarios.

B Offload Embedding Lookup to LMDB

In VocabTailor, model embeddings are offloaded to the CPU to reduce GPU memory usage. However, on devices with limited CPU memory or unified CPU-GPU memory architectures, this design alone does not sufficiently mitigate memory consumption. Such constraints are prevalent in edge devices, where SLMs are frequently deployed, as matrix multiplications can still be executed on the CPU within acceptable latency. In this context, VocabTailor’s core principle—decoupling and offloading embedding weights to a lower storage hierarchy—remains effective for two primary reasons: (1)

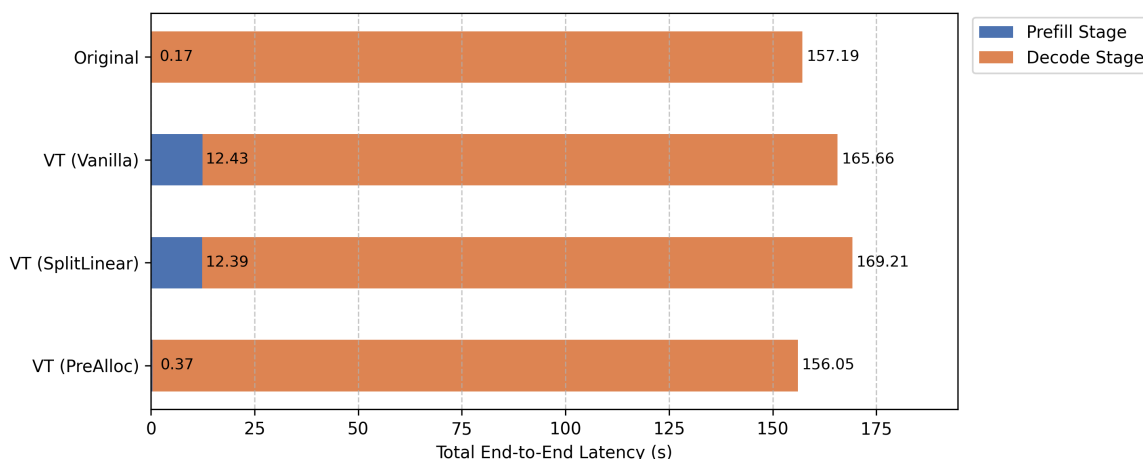


Figure 4: Latency decomposition across different dynamic LM head construction approaches. PreAlloc (bottom) successfully reduces the total prefill time from 12.43s (Vanilla) back to 0.37s, effectively matching the latency profile of the original model.

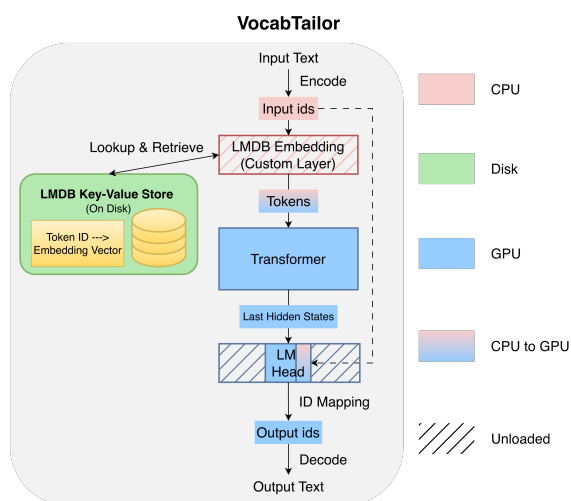


Figure 5: Overview of the VocabTailor framework with disk-backed embedding offloading. The embedding layer is replaced with a custom embedding layer that retrieves the corresponding tokens from the LMDB when the forward is called.

embedding lookup operations are computationally lightweight, and (2) the dominant inference latency arises from the forward passes, particularly tensor multiplications, making data movement a secondary bottleneck. Building on this principle and inspired by prior work on memory-mapped embedding storage (Yu et al., 2025), we implement disk-backed embedding offloading using the Lightning Memory-Mapped Database (LMDB). The embedding weights are serialized as key-value pairs, with token indices as keys and the corresponding embedding vectors as values. A disk-backed LMDB database is constructed to store these pairs, en-

abling on-demand retrieval of embeddings during inference. The extended framework is illustrated in Figure 5.

B.1 Experimental Setup

We evaluate the effect of embedding offloading for the machine translation (English-to-Chinese) task on a subset of 100 examples from the WMT24++ dataset (Deutsch et al., 2025). We use the Qwen3-0.6B (Yang et al., 2025) as the base model. All experiments are conducted on a single NVIDIA A100 GPU. Building on our previous results identifying PreAlloc as the most efficient dynamic LM head construction strategy, we compare three configurations that progressively incorporate embedding offloading: (1) Original, the original model with standard full-vocabulary; (2) VT (PreAlloc), which offloads embedding weights to CPU and employs a pre-warmed GPU buffer of size 128; and (3) VT (PreAlloc) + DiskEmb, which further offloads embeddings to disk-backed LMDB storage. Performance is measured across 100 prompts to capture statistical distributions of latency metrics including Time to First Token (TTFT), Time Per Output Token (TPOT), and end-to-end (E2E) latency. We also report throughput metrics including Token Per Second (TPS), Request Per Second (RPS), and average output length. Memory usage is monitored via the PyTorch CUDA memory management interface to report peak VRAM consumption.

B.2 Results

Table 8 presents a comparison between the three configurations. The latter two configurations share

the same dynamic LM head construction strategy and differ only in how embedding weights are stored and accessed.

Latency behavior. Compared to the Original, VT (PreAlloc) increases mean TTFT from 1.74 ms to 3.71 ms, reflecting the additional overhead of dynamic LM head construction and CPU-based embedding lookups. Introducing disk-based offloading further increases TTFT slightly to 4.58 ms, as embedding vectors are retrieved from disk rather than CPU memory. Importantly, this overhead remains small in absolute terms.

Tail latency analysis reveals a different trend. While VT (PreAlloc) exhibits a higher P99 TTFT (17.20 ms), DiskEmb configuration reduces P99 TTFT to 12.23 ms, indicating more stable prefill latency. This suggests that disk access introduces predictable overhead without amplifying variance across inputs.

TPOT remains comparable across all configurations, with mean TPOT values clustered around 24 ms. This confirms that embedding offloading—whether to CPU or LMDB—does not affect steady-state decoding, which is dominated by transformer forward passes.

End-to-end (E2E) latency remains broadly consistent across configurations. Despite differences in TTFT, both VT-based models achieve comparable mean and slightly lower tail E2E latency relative to the Original, indicating that prefill overhead is amortized over the generation process for typical output lengths.

Latency decomposition. Figure 6 illustrates the breakdown of total end-to-end (E2E) latency consumed by the prefill stage versus the decode stage. In the vanilla configuration, the prefill stage is nearly instantaneous, accounting for only 0.11% of total execution time. Transitioning to VT (PreAlloc) increases this proportion slightly to 0.24%. When moving further to DiskEmb, the prefill percentage still remains remarkably low at 0.30%. This demonstrates that while retrieving embedding vectors from disk-based storage (LMDB) is technically more time-consuming than CPU memory access, the impact on the overall latency profile is negligible. In all cases, the decoding stage continues to dominate over 99.7% of the runtime.

Throughput and resource efficiency. Table 9 reports decode throughput and resource usage. Decode TPS remains stable across configurations (around 42 tokens/s). This supports the observation that offloading embedding does not degrade decod-

ing performance. RPS is slightly lower for VT variants, consistent with their higher per-request initialization overhead. Peak GPU memory is reduced from 1.17 GB (Original) to 0.91 GB for both VT-based configurations, reflecting the benefit of pre-allocating LM head weights on GPU and offloading embeddings. VT (PreAlloc) has embedding weights on the CPU, which consumes 0.28 GB of memory. With DiskEmb, the embedding is evicted to disk, eliminating the CPU memory overhead.

Summary. Overall, the results confirm that embedding offloading is compatible with VocabTailor and provides additional memory savings. Both CPU-based and disk-backed offloading improve GPU memory efficiency relative to the original model without compromising decoding performance, making VocabTailor suitable for memory-constrained or resource-limited environments.

C Support Both Tied and Non-Tied Embedding Architectures

In large language models, the embedding layer and the LM head share the same weight dimension. The hidden states after going through the LM head and the input tokens before feeding into the transformer layers are considered in a similar representation space. Thus, to reduce the total number of parameters, some models would share the weight of the embedding and the LM head. Weight tying effectively reduces model size and inference-time memory consumption and is therefore widely adopted in small language models (SLMs). However, a portion of SLMs still retains non-tied embeddings to keep higher expressiveness.

VocabTailor supports both tied and non-tied embedding architectures, with their respective inference workflows illustrated in Figure 7. For tied embedding models, token lookup during the prefill stage is performed using the embedding weights stored on the CPU. During decoding, because the embedding and LM head weights are shared, we can directly use the reduced LM head weights to replace the embedding weights without extra memory usage, thereby eliminating memory movement overhead. Moreover, since the LM head is already reduced, embedding lookup can be performed without an ID mapping between the original token index and the reduced version token index. Such mapping is deferred until generation completes, before tokenizer decoding.

Model	TTFT (ms)				TPOT (ms)				E2E Latency (s)			
	Mean	P50	P90	P99	Mean	P50	P90	P99	Mean	P50	P90	P99
Original	1.74	0.79	0.90	2.07	24.21	23.64	24.09	25.17	1.57	1.56	2.44	4.12
VT (PreAlloc)	3.71	2.63	3.35	17.20	24.30	23.82	24.32	25.96	1.56	1.57	2.38	3.91
VT (PreAlloc) + DiskEmb	4.58	3.64	4.10	12.23	24.46	24.01	24.27	25.50	1.57	1.58	2.39	3.82

Table 8: Latency distribution for three inference configurations on a machine translation (English-to-Chinese) task ($N = 100$), comparing the effect of embedding offloading. We report Time to First Token (TTFT), Time Per Output Token (TPOT), and end-to-end (E2E) latency. Percentiles (Mean, P50, P90, P99) are included to characterize tail latency and model stability under varying prompt complexity.

Model	TPS (token/s)	RPS (req/s)	Avg. Output Length	Weights on CPU (GB)	Peak VRAM (GB)
Original	41.86	0.64	65.73	0	1.17
VT (PreAlloc)	41.65	0.55	64.84	0.28	0.91
VT (PreAlloc) + DiskEmb	41.44	0.51	64.78	0	0.91

Table 9: Throughput and resource utilization for three inference configurations on a machine translation (English-to-Chinese) task ($N = 100$), comparing the effect of embedding offloading. Decode throughput is measured in Tokens Per Second (TPS), computed as total generated tokens divided by total decoding time. Request throughput (RPS), average output length, weights on CPU, and peak GPU memory consumption are also reported.

For the non-tied embedding models, the embedding layer remains on the CPU and is used for both prefill and decode stages. Since only the LM head is reduced, the token indices produced by the reduced LM head after the softmax and sampling operation must be mapped back to the original indices before calling the embedding layer forward to get the corresponding token vectors. After generation, the produced token IDs are already the same as the ones in the full vocabulary, so no additional ID mapping is required for the tokenizer decoding.

D Profiling Strategy Algorithm

The detailed profiling strategy is presented in Algorithm 1.

E Discussion

E.1 Discussion on Memory Reduction

VocabTailor achieves substantial reductions in the memory usage of vocabulary-related components. In certain tasks (e.g., information extraction), this can reduce the active vocabulary to a small fraction of the full set, yielding up to 99% savings. However, these components (i.e., embedding matrix and LM head) account for only a portion of the total model memory, which also includes transformer layers, activations, and KV cache. As a result, the overall reduction in end-to-end VRAM usage is more moderate (22-23% as shown in Ta-

ble 2, though still meaningful in practice. This distinction highlights the importance of separating component-level savings from system-level memory usage when evaluating efficiency gains.

E.2 Discussion on Task Switching and Switching Cost

In practical small language model (SLM) deployment scenarios, models are typically optimized for a single downstream task or a small set of closely related tasks. High-frequency switching between substantially different task-specific vocabularies is relatively uncommon in such settings. In many production settings, SLMs are instantiated as task-specialized services, thereby avoiding repeated re-configuration during runtime. In scenarios where task switching is required, such switching is generally low-frequency. In these cases, the overhead associated with loading task-specific LM Head parameters remains modest relative to overall deployment and inference costs. For VocabTailor, it takes 0.2 s to switch between task-specific vocabularies, which we consider acceptable in low-frequency scenarios.

E.3 Discussion on Out-of-distribution

Out-of-distribution is a common issue for existing vocabulary pruning methods. In fact, one of our contributions is to mitigate this issue. Out-of-distribution is bounded by our design as VocabTai-

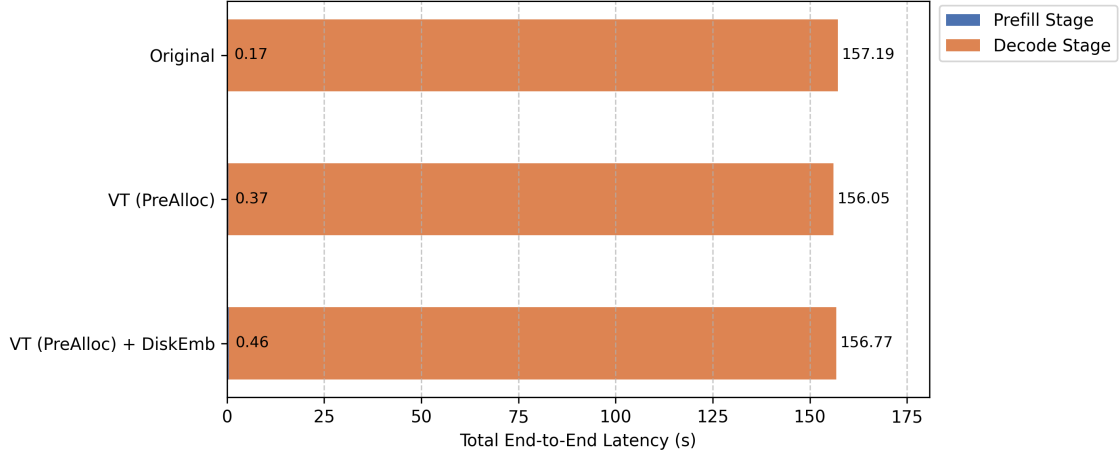


Figure 6: Latency decomposition of three inference configurations. We report the total end-to-end (E2E) latency consumed by the prefill (initialization) stage versus the decoding (generation) stage. While offloading embeddings to the CPU or LMDB increases the prefill proportion relative to the Original (0.11%), the impact remains marginal (0.24% for VT PreAlloc and 0.30% for DiskEmb), ensuring that the latency remains dominated by decoding.

Algorithm 1 Profiling Strategy

Definition:

- \mathcal{V} : Corpus-profiling vocabulary
- M : Total number of documents in the corpus
- \mathcal{I}_i : Set of input tokens corresponding to a single example i
- \mathcal{U} : Set of language-specific tokens obtained from Unicode blocks
- $df \in \mathbb{R}$: Dictionary mapping token v to its document frequency $df(v)$
- $\tau \in [0, 1]$: Tolerance threshold (fraction of documents allowed to be impacted)
- \mathcal{T} : Final calibrated task-specific vocabulary

Input: \mathcal{V}, τ

Output: \mathcal{T}

- 1: **Input-Aware Filtering**
 - 2: $\mathcal{V}_1 \leftarrow \{v \in \mathcal{V} \mid v \notin \mathcal{I}_i\}$
 - 3: **Language-specific Filtering**
 - 4: $\mathcal{V}_2 \leftarrow \{v \in \mathcal{V}_1 \mid v \in \mathcal{U}\}$
 - 5: **Tolerance Filtering**
 - 6: $N \leftarrow |\mathcal{V}_2|$
 - 7: $F \leftarrow$ list of $(v, df(v))$ for $v \in \mathcal{V}_2$
 - 8: Sort F ascending by $df(v)$
 - 9: $\text{index} \leftarrow 0$
 - 10: **for** $j \leftarrow 1$ **to** $N + 1$ **do**
 - 11: **if** $\text{sum}(F[:j]) > \tau M$ **then**
 - 12: $\text{index} \leftarrow N - j + 1$
 - 13: **break**
 - 14: **end if**
 - 15: **end for**
 - 16:
 - 17: **return** $\mathcal{T} \leftarrow \{v \in F[-\text{index} :]\}$
-

lor retains the full tokenizer. (details in Section 3.1), dynamically loads the input tokens for each example at runtime, and maintains a task-specific static vocabulary T . Thus, in the prefilling stage, the input text is encoded as the same tokens and generates identical hidden states as the original model without vocabulary pruning. In essence, our design eliminates the out-of-distribution problem for the prefilling stage (which happened for the current static vocab pruning). The only potential limitation concerns rare output-only tokens not captured in T . However, tolerance filtering (Section 3.4) explicitly bounds worst-case degradation via τ . In our experiments ($\tau = 0.01$), no catastrophic failure was observed across five diverse tasks (Table 1).

F Data Processing and Evaluation Details

For machine translation, we follow the best practice of zero-shot machine translation by Zhang et al. (2023), using a simple English template.

English-to-Chinese translation template:

```
[{"role": "user", "content": "English:\n{SOURCE}\n_Chinese:"}]
```

English-to-Italian translation template:

```
[{"role": "user", "content": "English:\n{SOURCE}\n_Italian:"}]
```

For summarization, the model is finetuned and evaluated using prompt shown below.

Summarization prompt template:

```
[{"role": "user", "content": "f\"Document:\n{DOCUMENT}\n"}]
```

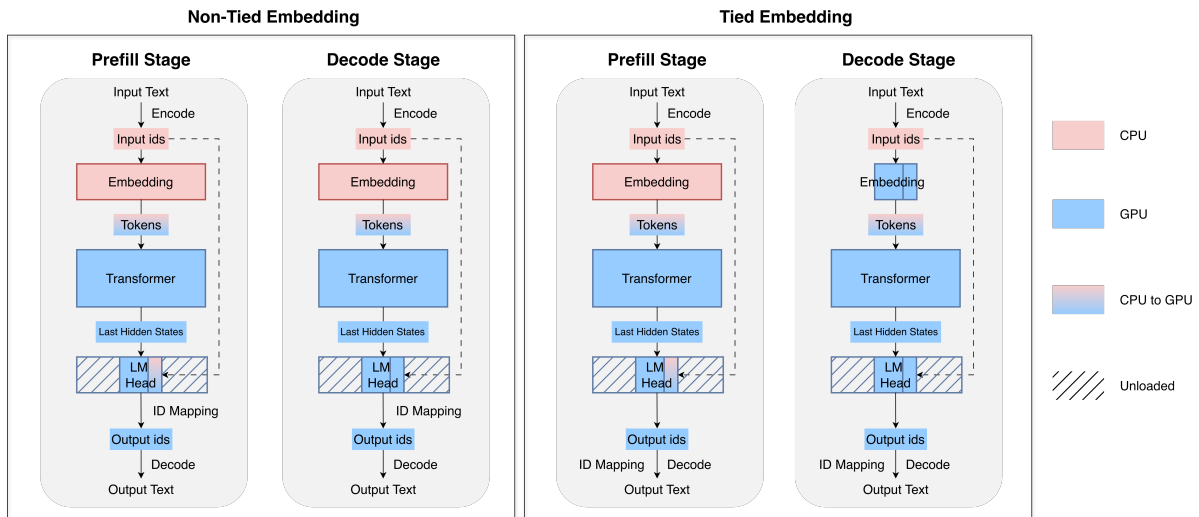


Figure 7: Comparison of VocabTailor workflows. Left: In tied architectures, LM head weights are reused for decoding to eliminate memory movement. ID mapping is deferred until the final tokenizer stage. Right: In non-tied architectures, only the LM head is reduced, so an explicit ID mapping is required after sampling to align the reduced LM head and the full CPU embedding layer.

For the code completion task, we use the evaluation script of SAFIM. For information extraction, we use the LM-Eval-Harness framework and set the task to ‘squad_completion’. For math word problem solving, we use the Math-Eval-Harness framework and set the task to ‘mawps’.

G Fine-tuning Config for Summarization Baseline Model

For summarization, all the methods are based on a finetuned Llama 3.2 3B model. We set learning rate=2e-5, batch size = 128. The model is finetuned using XSUM train split set for 1 epoch.

H Use of AI for Writing/Coding Assistance

Generative AI use in this work is limited to assistance purely with the language of the paper and short-form input assistance under the AI writing assistance policy.