
zip2zip: Inference-Time Adaptive Vocabularies for Language Models via Token Compression

Anonymous Authors¹

Abstract

Tokenization efficiency plays a critical role in the performance and cost of large language models (LLMs), yet most models rely on static tokenizers optimized on general-purpose corpora. These tokenizers' fixed vocabularies often fail to adapt to domain- or language-specific inputs, leading to longer token sequences and higher computational costs. We introduce zip2zip, a framework that enables LLMs to dynamically adjust the token vocabulary at inference time, allowing for fewer generated tokens and thus faster inference. zip2zip consists of three key components: (1) a tokenizer based on Lempel-Ziv-Welch (LZW) compression that incrementally merges co-occurring tokens into reusable *hypertokens* on the fly; (2) an embedding layer that computes embeddings for newly formed hypertokens at runtime; and (3) a causal language modeling variant that trains the model to operate on hypertokenized, compressed sequences. We show that an existing LLM can be zip2zip-fied in 10 GPU-hours via parameter-efficient finetuning. The resulting zip2zip LLMs effectively learn to use hypertokens at inference time, reducing input and output sequence length by 20–60%, with significant improvements in inference latency. Code is available at this [anonymous repository](#).

1. Introduction

Large language models (LLMs) have shown impressive versatility across a broad spectrum of tasks and domains (Brown et al., 2020; Bubeck et al., 2023), including biomedical tests (Nori et al., 2023), mathematical reasoning (Frieder et al., 2023), programming (Jiang et al., 2024), and multiple human languages. A critical underlying com-

ponent of this flexibility is the tokenizer, which defines the model's vocabulary and governs how raw text is converted into token sequence fed to the model. The efficiency of the tokenization scheme—i.e., how compactly a text is represented as tokens—has significant impact on model performance. In particular, a more *compact* tokenization yields three key benefits: (1) larger effective context window; (2) lower computational (and thus monetary) cost; and (3) shorter response times.

Despite its importance, the tokenizer used in most LLMs produces a fixed, static vocabulary using algorithms such as Byte Pair Encoding (Sennrich et al., 2016) over large-scale, general-purpose web corpora. While this globally optimized vocabulary performs reasonably well on average, it often fails to adapt to domain-specific or language-specific distributions (Ahia et al., 2023; Petrov et al., 2023), where the text distribution diverges significantly from the pretraining data. The resulting mismatch leads to longer token sequences, increasing both memory and compute demands, as well as the end user's cost by a factor of 2–3x when processing domain-specific text (Ahia et al., 2023). To mitigate this issue, prior work has explored expanding the token vocabulary during domain or language adaptation to improve tokenization efficiency (Wang et al., 2019; Zhao et al., 2024; Kim et al., 2024; Liu et al., 2023; 2024). While effective, this approach needs to be repeated for each target domain or language and requires maintaining separate tokenizers. Meanwhile, commercial LLM providers trend toward increasing the size of token vocabularies—growing from 32K to 128K (Grattafiori et al., 2024) and even up to 200K (Abdin et al., 2024) tokens—to improve overall tokenization efficiency. However, prior work (Dagan et al., 2024; Liang et al., 2023) shows that simply enlarging the vocabulary yields diminishing returns in domain adaptation, and vocabularies past a certain size can potentially degrade model performance (Liang et al., 2023).

These limitations point to a compelling need for an adaptive tokenization mechanism—one that can dynamically tailor the vocabulary to the input text at inference time, without retraining the model or maintaining separate tokenizers. Such a mechanism would allow the model to construct new domain-specific tokens on-the-fly, so to enhance tokeniza-

¹Anonymous Institution, Anonymous City, Anonymous Region, Anonymous Country. Correspondence to: Anonymous Author <anon.email@domain.com>.

Preliminary work. Under review by the International Conference on Machine Learning (ICML). Do not distribute.

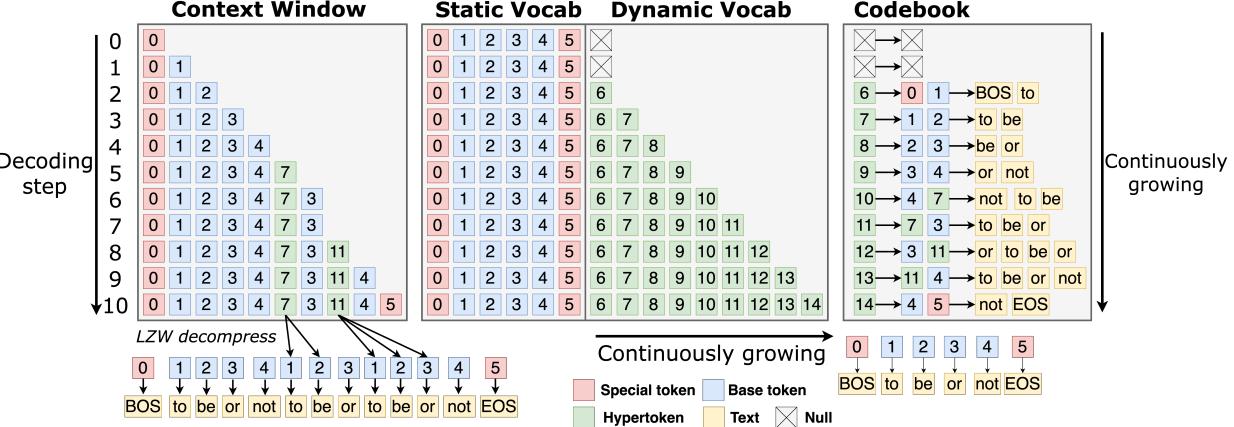


Figure 1: Overview of the zip2zip inference pipeline. At each decoding step, the model has a growing **context** composed of both **base tokens** (blue) and **hypertokens** (green). The **static vocabulary** of size 6 remains fixed, while the **dynamic vocabulary** is continuously expanded by merging co-occurring tokens using **LZW compression**. The **codebook** (right) maps hypertoken IDs to their corresponding base tokens. As decoding progresses, new hypertokens created at step t (e.g., “to be”, “or not”) become **immediately** available for reuse at step $t+1$. Additionally, output tokens, once generated, **instantly** become eligible for compression. Hypertokens are also eligible for merging, enabling the formation of **nested hypertokens**. The final output sequence (bottom) is reconstructed via LZW decompression.

tion efficiency. However, adaptive tokenization poses architectural challenges, as both the embedding layer and the language modeling head in transformer models (Vaswani et al., 2017) are static matrices tied to a fixed vocabulary size.

In this paper, we propose zip2zip (with a hat-tip to seq2seq (Sutskever et al., 2014)), a method that equips LLMs with a dynamic token vocabulary, enabling inference-time token adaptation. zip2zip achieves adaptive tokenization through *continuous vocabulary expansion* at runtime, allowing the model to represent a repeated or domain-specific pattern with a single long token rather than inefficient short tokens. This requires modest modifications to both the transformer architecture and the language modeling objective. zip2zip comprises three key components: (1) **Tokenizer**: an integration of Lempel-Ziv-Welch (LZW) compression¹ (Welch, 1984) into the tokenization process, which continuously merges frequently co-occurring token sequences into reusable longer tokens (hypertokens) at runtime; (2) **Architecture**: a lightweight encoder added to the transformer that computes embeddings for newly formed tokens on the fly; (3) **Training**: a compression-aware causal language modeling variant that trains the model directly on LZW-compressed sequences, aligning learning with the inference-time token distribution. The name zip2zip reflects its dual role in achieving compression of both the input tokens (the first *zip*) and output tokens (the second *zip*), thereby jointly improving the efficiency of input encoding and output decoding. We finetune Phi-3-4B and Phi-3-14B to support zip2zip using as few as 100M

¹LZW is the algorithm used in zip compression tool, which inspired the name zip2zip.

tokens—requiring only 10 and 40 H100 GPU hours, respectively—for effective adaptation. The resulting models demonstrate strong inference-time compression capabilities and achieve 20–60% reductions in both input and output sequence lengths, translating to up to 60% improvements in end-to-end latency.

To make it easy to upgrade existing LLMs to zip2zip, we release an efficient, open-source implementation of the framework. It includes (1) a fast Rust-based LZW tokenizer, (2) a drop-in model architecture compatible with Hugging Face Transformers and vLLM, (3) a training pipeline for LZW-compression-based finetuning. Existing LLMs can be seamlessly extended with zip2zip, gaining adaptive tokenization capabilities through parameter-efficient finetuning—without any changes to the base model or tokenizer.

2. zip2zip

2.1. Dynamic Token Vocabulary

To enable dynamic tokenization at inference time, we associate each LLM with a *hyper-vocabulary* \mathcal{V}_h that augments the model’s static token vocabulary. Tokens from the original vocabulary \mathcal{V} are referred to as *base tokens*. Each entry in the hyper-vocabulary is a *hypertoken*, representing a merged sequence of base tokens. The total vocabulary for a zip2zip model is the union $\mathcal{V} \cup \mathcal{V}_h$. At the beginning of each inference session, \mathcal{V}_h is initialized as an empty set, and is incrementally populated during decoding by identifying and merging recurring token subsequences in the context window, as illustrated in Figure 1.

110
111
112
113
114
115
116
117
118
119
120
121
122
123
124
125
126
127
128
129
130
131
132
133
134
135
136
137
138
139
140
141
142
143

Continuous Vocabulary Expansion. As decoding proceeds, zip2zip continuously merge co-occurring tokens as new hypertokens to \mathcal{V}_h and recurrently apply merging on newly generated tokens. This *continual expansion* allows the model to represent longer, recurring sequences of base tokens compactly. Hypertokens are treated as first-class tokens within the model, used interchangeably with base tokens throughout the decoding process. Importantly, this process occurs entirely during inference, without modifying the underlying tokenizer or requiring model retraining.

LZW Algorithm. We implement vocabulary expansion using the Lempel-Ziv-Welch (LZW) compression algorithm—a dictionary-based, lossless compression method that incrementally builds a codebook of variable-length sequences. In our setting, the codebook is initialized with the base token vocabulary \mathcal{V} and expands by adding new hypertokens on the fly as recurring token patterns are encountered. To control the growth of the dynamically expanding vocabulary, we impose a maximum merge size M that restricts how many base tokens a single hypertoken can represent. LZW is particularly well-suited for zip2zip due to the following properties:

(1) it is *streaming*—hypertokens created at step t can be immediately reusable at step $t + 1$; in contrast, methods like BPE require access to the full sequence and operate offline;

(2) it is *self-contained*—input base tokens can be perfectly reconstructed from the compressed token sequence alone²;

(3) it is *unambiguous*—when both base tokens and hypertokens are available, which one to use is consistently determined by the LZW algorithm without ambiguity.

2.2. Hyper-Embedding and Hyper-Projection

Hypertokens do not have fixed embedding vectors in the original model’s embedding layer (and projection layer), as they are not part of the original vocabulary. To compute the embedding of a hypertoken, we learn a mapping from the base token embeddings to the hypertoken embedding. We achieve this by introducing a *hyper-encoder*, which is a neural network that takes the embeddings of the constituent base tokens as input and outputs the corresponding hypertoken embedding. Specifically, for a sequence of M base tokens $y_{1:M} := y_1 \dots y_M$, the hyper-encoder $f_\phi : \mathcal{V}^M \rightarrow \mathbb{R}^d$ produces the hypertoken embedding $h = f_\phi(y_{1:M}) \in \mathbb{R}^d$, where M is the maximum merge size and d is the embedding dimension. For hypertokens composed of fewer than M base tokens, we pad the input sequence to length M . Since the embedding map for base tokens remains unchanged, the hyper-encoder f_ϕ essentially maps the concatenated base

²There is no need to persist or transmit the codebook across inference calls, preserving compatibility with existing LLM libraries and interfaces.

token embeddings from a $(M \times d)$ -dimensional space to a d -dimensional hypertoken embedding vector, performing nonlinear dimensionality reduction.

For the output projection layer, if the underlying transformer ties the embedding and the projection matrices, one can reuse the same hyper-encoder to compute the representation used for projection. Otherwise, a separate hyper-encoder is trained to produce the hypertoken projection vectors.

2.3. Architecture

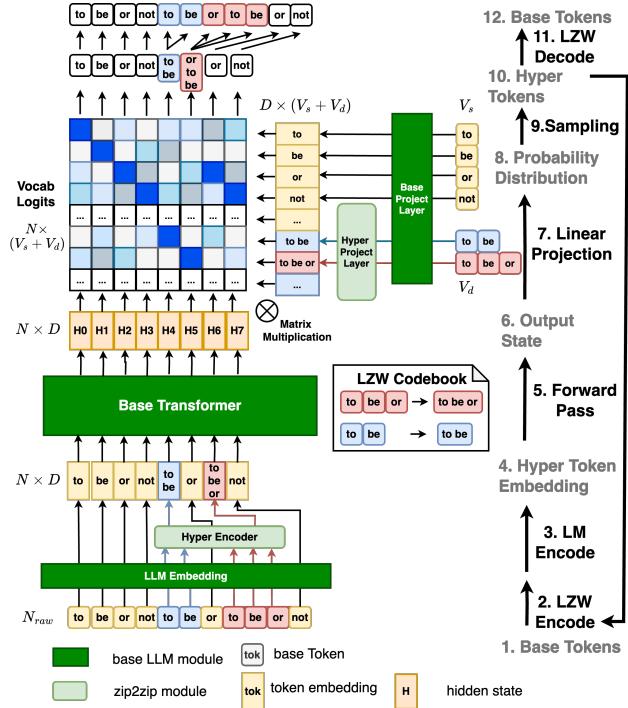


Figure 2: zip2zip architecture. At inference time, base tokens are compressed into hypertokens using LZW (**STEPS 1–2**). A hyper-encoder computes embeddings for hypertokens (**STEP 3–4**), which are processed by the base LLM (**STEPS 5–6**). Output representations are projected jointly on base and hyper-projection layers (**STEP 7**), producing joint logits and sampled tokens (**STEPS 8–10**), which can be decoded back to base tokens (**STEPS 11–12**).

We illustrate the architecture of zip2zip in Figure 2. The input text is first tokenized into base tokens (**STEP 1**), which are then passed through an online LZW compressing module that compresses the token sequence into a stream of hypertokens (**STEP 2**). Since hypertokens are not part of the model’s original embedding layer, their embedding vectors are computed on-the-fly using the *hyper-encoder* during inference (**STEP 3–4**). Once embedded, both base token embeddings and hypertokens embeddings are passed through the standard transformer layers of the base model, producing contextualized hidden states (**STEP 5–6**). This step is identical to vanilla transformer, with hypertokens and base tokens treated equally. At the output projection layer, hyper-

165 token projection vectors (same as the hypertoken embedding
 166 vectors in the tied case, and computed by a separate hyper-
 167 encoder otherwise) are appended to the original projection
 168 matrix in the language modeling head (**STEP 7**). This allows
 169 the model to compute a joint softmax over the union of the
 170 base vocabulary and the hyper vocabulary $\mathcal{V} \cup \mathcal{V}_h$ (**STEP**
 171 **8**). The resulting probability distribution is over $\mathcal{V} \cup \mathcal{V}_h$,
 172 and the sampled token may be either a base token or a hyper-
 173 token (**STEP 9**). In the next cycle, the newly generated
 174 token (**STEP 10**)—whether base or hyper—is appended to
 175 the input sequence, and the process repeats (back to **STEP**
 176 **1**). At the end of generation, the hypertoken sequence is
 177 decompressed via the LZW decoding function into a se-
 178 quence of base tokens (**STEP 11–12**). The whole process
 179 works in a fully *autoregressive* way, where newly generated
 180 tokens will also be merged into hypertokens for future steps.
 181 Furthermore, we highlight two points:

182 **Consistent Vocabulary Updates.** The expanding vocab-
 183 uary—comprising newly created hypertokens—must be
 184 updated in a *consistent* manner across both the input em-
 185 bedding layer and the output projection layer, maintaining
 186 a consistent view of the hypertoken set. Failure to update
 187 both sides consistently can result in two types of errors:
 188 (1) hypertokens that cannot be decoded, or (2) the model
 189 attempting to decode a non-existing hypertoken.

190 **Hyper-Embedding Cache.** Although hypertoken embed-
 191 dings are computed on-the-fly, they are context-independent
 192 and can thus be cached across inference steps. Similar to
 193 the transformer’s KV-cache, this enables *incremental* up-
 194 dates: only newly created hypertokens need to be processed
 195 at each step. Since the codebook grows linearly with the
 196 number of tokens in the context n , the total cache size grows
 197 also linearly in memory. Thus, the computational cost for
 198 hypertoken embeddings remains constant per step—i.e., one
 199 token embedding is computed per step.

200 201 202 2.4. Training zip2zip models

203 **Objective.** Let \mathcal{D} denote the target text distribution. Given a
 204 language model π_θ parameterized by θ , standard pretraining
 205 seeks to minimize the causal language modeling (CLM)
 206 objective, which corresponds to the expected negative log-
 207 probability of data sequences under the model:

$$208 \min_{\theta} \mathbb{E}_{y \sim \mathcal{D}} [-\log \pi_\theta(y)], \quad (1)$$

211 where $\pi_\theta(y)$ denotes the probability of the token sequence
 212 y under the model π_θ .

213 Let \mathcal{C} be an *online* compression algorithm (e.g., LZW),
 214 and ϕ be the parameters of the hyper-encoder. Given a
 215 sequence $y \sim \mathcal{D}$, let $z = \mathcal{C}(y)$ be its compressed form.
 216 In zip2zip, we aim to optimize the same CLM loss, but
 217 over the compressed sequences z . The training objective

218 becomes:

$$\begin{aligned} & \min_{\theta, \phi} \mathbb{E}_{y \sim \mathcal{D}} [-\log \pi_{\theta, \phi}(\mathcal{C}(y))] \\ &= \min_{\theta, \phi} \mathbb{E}_{z \sim \mathcal{C}(\mathcal{D})} [-\log \pi_{\theta, \phi}(z)]. \end{aligned} \quad (2)$$

219 Here, we slightly abuse the notation to let $\pi_{\theta, \phi}(z)$ denote the
 220 probability assigned to the compressed sequence z , parame-
 221 terized by the base model weights θ and the hyper-encoder
 222 parameters ϕ .

223 To construct the compressed dataset $\mathcal{C}(\mathcal{D})$, we first tokenize
 224 the corpus using a standard tokenizer, and then apply the
 225 LZW compression algorithm. This preprocessing step is
 226 performed once prior to training and can be efficiently par-
 227 allelized through batching.

228 **Parallelizable Training via Causal Masking.** Although hy-
 229 pertokens introduce additional vocabulary dynamics, train-
 230 ing remains fully parallelizable. We leverage the standard
 231 causal masking mechanism used in language models, allow-
 232 ing the model to predict the next token—whether a base
 233 token or a hypertoken—at each position in parallel. To
 234 eliminate the need for sequential codebook updates during
 235 inference, we precompute a fixed codebook by applying
 236 LZW compression to the entire input sequence. This pre-
 237 computed codebook is then used consistently throughout
 238 training to condition token predictions, ensuring efficiency
 239 and compatibility with standard training pipelines.

240 **Auxiliary Reconstruction Loss.** We introduce an auxiliary
 241 reconstruction objective that encourages a hypertoken em-
 242 bedding to retain sufficient information about its underlying
 243 base token sequence. Specifically, the model is trained to
 244 reconstruct the original base token embeddings from the
 245 hypertoken embedding. We jointly optimize the language
 246 model and the hyper-encoder using a combined loss that
 247 includes both the standard next-token prediction loss and
 248 the auxiliary reconstruction loss. Formally, we optimize:

$$\begin{aligned} & \min_{\theta, \phi, \psi} \mathbb{E}_{y \sim \mathcal{D}} [-\log \pi_{\theta, \phi}(\mathcal{C}(y))] \\ &+ \lambda \mathbb{E}_{y_{1:M}} [\Delta(y_{1:M}, f_\psi(f_\phi(y_{1:M})))], \end{aligned} \quad (3)$$

249 where $f_\phi : \mathcal{V}^M \rightarrow \mathbb{R}^d$ is the hyper-encoder, $f_\psi : \mathbb{R}^d \rightarrow \mathcal{V}^M$
 250 is the decoder aiming to reconstruct the corresponding base
 251 tokens from their hyper-embedding, and $\Delta : \mathcal{V}^M \times \mathcal{V}^M \rightarrow$
 252 \mathbb{R} is the reconstruction loss function, such as the cross-
 253 entropy loss, between the base tokens $y_{1:M}$ and the recon-
 254 structed base tokens $f_\psi(f_\phi(y_{1:M}))$. The hyperparameter
 255 λ controls the trade-off between the prediction error of the
 256 language model and the reconstruction error of the autoen-
 257 coder. This joint optimization objective encourages the
 258 hyper-encoder to learn a compact d -dimensional manifold
 259 embedded in the higher-dimensional $(M \times d)$ space of base
 260 token embeddings, while the language model $\pi_{\theta, \phi}$ learns to
 261 predict the next (hyper)token given the preceding context.

	<Python code>	<Biomedical>	<French>
220 221 222 223 224 225 226 227 228 229 230 231 232	<pre>Implement a Transformer model in Pytorch that takes a sequence of tokens and outputs a hidden representation for each token. The Transformer should consist of an embedding layer, a self-attention mechanism, and a feed-forward neural network. Ensure that the Transformer can handle variable input and output sizes, as well as different sequence lengths.</lend><assistant> Here is an example implementation of a Transformer model in PyTorch:</pre> <pre>... import torch from torch.nn import Linear, LayerNorm, Module, Unfold, ModuleList class TransformerLayer(Module): def __init__(self, d_model, nhead, dim_feedforward): super().__init__() self.d_model = d_model self.nhead = nhead self.dim_feedforward = dim_feedforward self.dropout_attn = nn.Dropout() self.attn = MultiheadAttn(d_model, nhead, dim_feedforward, dropout=dropout_attn) self.layernorm = LayerNorm(d_model) self.ffn = FeedForward(d_model, dim_feedforward) # TODO - implement the feedforward neural network</pre>	<pre><user> Please explain what is messenger ribonucleic acid</lend><assistant> Messenger RNA (mRNA) is a type of ribonucleic acid (RNA) that serves as the messenger carrying genetic information from DNA to the ribosomes, where the protein synthesis takes place. The genetic code on the DNA sequence is transcribed into mRNA in a process called transcription.</pre> <pre>As it moves out of the nucleus of the cell, the mRNA strand encounters the ribosomes that read the sequence and use it as a guide to create a protein in a process known as translation. mRNA has a sequence of nucleic acids according to which the cell makes the enzymes, vitamins, and many other enzymes needed for cellular function. Therefore, it is vital in all biological functions since proteins are required by the cell.</pre> <pre>To summarize, mRNA is transcribed from DNA and is a type of RNA that serves as a messenger to carry genetic information within the cell. It instructs ribosomes on how to build particular proteins and, in turn, carries out the process of protein synthesis.</lend><user> Can you give me some examples of how to identify mRNA in laboratory experiments?</lend><assistant> Sure! Here are some examples of how to identify mRNA in laboratory experiments:</pre>	<pre><user> Expliquez-moi l'histoire de la Tour Eiffel</lend><assistant> La Tour Eiffel est une tour à étages émaillée située à Paris, France. Elle a été conçue par l'ingénieur Gustave Eiffel, dont le nom est associé à la construction de cette tour. La construction de la tour a été achevée en 1889 pour célébrer le 100e anniversaire de la Révolution française.</pre> <pre>La Tour Eiffel, haute de 324 mètres, est devenue le symbole de Paris et son point de repère le plus célèbre dans le monde. Elle est également l'un des monuments les plus visités au monde avec plus de 7 millions de visiteurs chaque année. La tour est utilisée comme balcon panoramique offrant des vues magnifiques sur Paris et la région.</pre> <pre>Malgré les critiques initiales de nombreux Parisiens sur la conception proéminente de cette tour, elle reste le plus populaire du monde. En effet, la Tour Eiffel est la seule structure du type pyramidal avant l'escalier mobile de visiter ses côtés.</pre> <pre>La Tour Eiffel a subi de nombreux changements et modifications au cours de ses 133 ans de vie. En 1981, elle a reçu la distinction "Monument Historique" par le ministre de la Culture France. Et après la Seconde Guerre Mondiale, la tour est restée la carte par l'attracteur touristique de la ville de Parisiens.</pre>

Figure 3: Examples of zip2zip model outputs. The green text indicates hypertokens and the red text indicates the base tokens with magnified annotations illustrating hypertoken formation.

The reconstruction loss can be viewed as a form of auto-encoding, where the hypertoken acts as a compressed latent representation and reconstruction encourages the preservation of semantic content and the compression to be lossless.

Adapting Pretrained Language Models. The proposed objectives (Equation 2, 3) integrate naturally with pretrained language models. In this setting, the base model can be frozen while training only the hyper-encoder to adapt to compressed token sequences. Parameter-efficient methods such as LoRA (Hu et al., 2022) may also be used to adapt select components of the base model, enabling effective adaption with minimal computes.

2.5. Efficiency Advantage

zip2zip improves efficiency by increasing the average token length, thereby reducing the number of tokens required to represent the same text. This compression applies to both inputs (e.g., prompts) and outputs (e.g., completions), leading to shorter effective context lengths. As a result, the model performs fewer computations—both in the attention mechanism and the feedforward layers—and, more importantly, requires fewer autoregressive decoding steps during inference. Since the latency of large language models is primarily driven by the cost of sequential decoding, reducing the number of output tokens by $n\%$ leads to an approximate $n\%$ speedup in decoding latency, which we will demonstrate empirically in Section 3.6. A more detailed discussion of FLOPs is provided in Appendix B for completeness.

3. Experiments

To evaluate the effectiveness of zip2zip, we adapt the Phi-3 models (3B and 14B) within the zip2zip framework. We evaluate our adapted models across four dimensions: (1) token efficiency, (2) language modeling perplexity, (3) downstream task performance, (4) inference effi-

cency. For perplexity and downstream benchmarks, we use the widely adopted lm-evaluation-harness framework (Gao et al., 2024).

3.1. Training Setup

Rather than updating the full model weights, we adopt parameter-efficient finetuning using LoRA (Hu et al., 2022). In addition, we train the *hyper-embedding* and *hyper-projection* modules. We set the maximum merge size to $M = 3$ and use a two-layer transformer encoder as the hyper-encoder. Ablation studies on M and hyper-encoder architecture can be found in Appendix A. For comparison, we also perform continual finetuning of the base model using LoRA under identical training conditions, serving as a baseline (denoted as Cont. Finetune in the Tables). The finetuning process is highly efficient, requiring approximately 10 H100-GPU hours for a 4B-parameter model and up to 40 H100-GPU hours for a 14B-parameter model, using only 0.1 billion training tokens. Interestingly, the *reconstruction loss* converges to **near zero** during training, indicating that the model can almost perfectly recover the original base token sequences from the hypertoken representations. This highlights the learned compression is highly information-preserving. Details of the training setup, compute infrastructure, and dataset curation are provided in Appendices D and E.

3.2. Sample Outputs and Hypertoken Patterns

We present several examples to provide intuition into how the zip2zip model generates text.

We see that the model successfully generates a mixture of hypertokens and base tokens in the output (see Figure 3). The hypertoken ratio is as high as 40% in the Python code generation example, and 20% in the biomedical text generation example. Many of the hypertokens correspond to semantically meaningful units or domain-specific terms as

Table 1: Examples of hypertokens formed by zip2zip across three domains

Code Generation	Biomedical	French
Py + Torch = PyTorch	m + R + NA = mRNA	tour + istique = touristique
Layer + Norm = LayerNorm	mess + enger = messenger	le + plus = le plus
Trans + former = Transformer	gen + etic = genetic	Tour + Eiff + el = Tour Eiffel
d_ + model = d_model	D + NA = DNA	Paris + iens = Parisiens

Table 2: Token efficiency (bytes/token) across domains for different tokenizers w/wo zip2zip.

Tokenizer	Code	Math	Chat	Multilingual	Web
Llama-3-128K (Grattafiori et al., 2024)	4.1	2.7	5.1	3.8	4.6
+zip2zip	6.3 (+54%)	4.0 (+48%)	6.4 (+25%)	4.7 (+24%)	5.4 (+17%)
Qwen-2-150K (Yang et al., 2024)	4.0	2.3	5.1	3.7	4.4
+zip2zip	6.2 (+55%)	3.7 (+61%)	6.4 (+25%)	4.6 (+24%)	5.2 (+18%)
Phi-4-200K (Abdin et al., 2024)	4.1	2.7	5.4	4.6	4.7
+zip2zip	6.3 (+54%)	4.1 (+52%)	6.7 (+24%)	5.5 (+20%)	5.4 (+15%)
Gemma-3-256K (Team et al., 2025)	3.3	2.3	5.0	4.4	4.5
+zip2zip	5.6 (+70%)	3.7 (+61%)	6.4 (+28%)	5.4 (+23%)	5.4 (+20%)

shown in Table 1. For a more fine-grained visualization of hypertoken with zip2zip, we provide visualizations of token streams in Appendix 8.

3.3. Token Efficiency

Given an input text x and a tokenizer, we define *token efficiency* $\eta := \frac{\text{Bytes}(x)}{\text{Tokens}(x)}$ as the average number of bytes represented by each token, where Bytes(x) refers to the number of bytes in the UTF-8 encoding of x . This measures how compactly a tokenizer encodes input text—higher values of η indicate more efficient tokenization.

We evaluate token efficiency using the tokenizers of four LLMs—Llama-3 (Grattafiori et al., 2024), Qwen-2 (Yang et al., 2024), Phi-4 (Abdin et al., 2024), and Gemma-3 (Team et al., 2025)—each associated with a different base vocabulary size ranging from 128K to 256K. Token efficiency is measured across five representative domains, sampled from publicly available datasets: code (Lozhkov et al., 2024b), math (LI et al., 2024), chat (Ding et al., 2023), multilingual (Penedo et al., 2024), and web (Lozhkov et al., 2024a). Table 2 shows that applying LZW zip2zip consistently improves token efficiency across all tokenizer and domains. Gains are particularly strong in structured domains like code and math—50% higher than the base tokenizer. Interestingly, models with larger vocabulary sizes do not always achieve better token efficiency, suggesting that simply enlarging the vocabulary size is not sufficient to improve it.

3.4. Perplexity

We evaluate the perplexity of zip2zip models on four corpora: Wikitext (Merity et al., 2016), the Pile (Gao et al., 2020), and two subsets of Paloma (Magnusson et al., 2023): mC4, a multilingual subset of C4, and dC4 (aka C4-100D),

Table 3: Byte-perplexity (\downarrow) on four corpora using a 1024-token context window.

Model	Method	Wiki	Pile	mC4	DC4
Phi-3.5-4B	Base	1.62	1.88	1.94	1.77
	Cont. finetune	1.63	1.89	1.94	1.77
	zip2zip	1.71	2.02	2.04	1.84
Phi-3-14B	Base	1.43	1.72	1.82	1.67
	Cont. finetune	1.47	1.79	1.86	1.68
	zip2zip	1.56	1.90	1.96	1.75

a subset of C4 spanning 100 domains. Given a token sequence $x = x_1, \dots, x_N$, and a model q , perplexity and byte-level perplexity (Radford et al., 2019; Magnusson et al., 2023) are defined as: $\text{PPL} := \left(\prod_{i=1}^N q(x_i) \right)^{-1/N}$, $\text{Byte-PPL} := \left(\prod_{i=1}^N q(x_i) \right)^{-1/B} = \text{PPL}^{1/\eta}$, where B is the number of UTF-8 bytes of the text, and η denotes the token efficiency (i.e., bytes per token). Token-level perplexity depends on the tokenization scheme and is unsuitable for cross-tokenizer comparison. We instead report byte-level perplexity, a vocabulary-agnostic metric that normalizes for tokenization differences. Table 3 shows that zip2zip model has a modest increase in Byte-perplexity, indicating a drop in language modeling performance.

3.5. Evaluation on NLP Benchmarks

We next evaluate zip2zip’s performance on real-world tasks. We evaluate on seven widely used NLP benchmarks, including ARC-[Challenge, Easy] (Clark et al., 2018), HellaSwag (Zellers et al., 2019), LAMBADA (Paperno et al., 2016), OpenbookQA (Mihaylov et al., 2018), PIQA (Bisk et al., 2019), Winogrande (Sakaguchi et al., 2019) and

	Benchmark	Phi-3.5-4B			Phi-3-14B		
		Base	C.F.	Z2Z	Base	C.F.	Z2Z
ARC-c	0.60	0.60	0.57	0.62	0.62	0.62	0.62
ARC-e	0.83	0.82	0.83	0.80	0.88	0.86	
HS	0.66	0.63	0.61	0.70	0.66	0.68	
OBQA	0.46	0.47	0.46	0.51	0.52	0.51	
PIQA	0.79	0.82	0.82	0.83	0.87	0.85	
WG	0.75	0.75	0.75	0.76	0.80	0.79	
GSM8K	0.82	0.40	0.15	0.84	0.52	0.25	

Table 4: Two-shot accuracy (in %) across 7 NLP benchmarks. Higher is better. Standard deviations (bootstrapped) ≈ 0.02 across all tasks. C.F.=Continuous finetune, Z2Z=zip2zip.

GSM8K (Cobbe et al., 2021). As shown in Table ??, the model finetuned with zip2zip performs similarly to the baseline on most tasks. However, on GSM8K, where the primary task involves numerical computation, the model exhibits significant degradation. Due to the sensitivity of such tasks to tokenization, it occasionally generates malformed or repeated numbers. While token-level operations are already known to be challenging for LLMs (Singh & Strouse, 2024), adaptive tokenization appears to exacerbate this issue.

To validate the effectiveness of zip2zip on non-English languages, we evaluate the model on machine translation tasks, including WMT14 (Macháček & Bojar, 2014), WMT16 (Bojar et al., 2016). The results, shown in Table 5, indicate a small performance degradation across BLEU, CHRF, and TER metrics when using zip2zip. However, the drop is relatively minor, suggesting that the model retains strong multilingual capabilities even in the compressed representation.

3.6. Inference Efficiency

zip2zip reduces decoding time by lowering the number of tokens that need to be generated. However, it introduces additional FLOPs due to the on-the-fly computation of hyper-embeddings by the hyper-encoder. To address this overhead, we implement hyper-embedding caching and optimize the computation using a custom Triton kernel. We report separate timings for *prefill* and *decoding* across multiple models, with and without zip2zip, in Table 6.

As we show in Table 6, zip2zip achieves a significant speedup in all four settings. Both prefill and decoding times are significantly reduced, with the most substantial gains observed in the 512+256 setting with the Phi-3.5-4B model. Improvements are significantly stronger on datacenter-grade GPUs like the NVIDIA H100 and more modest on consumer hardware (e.g., Apple M1).

Table 5: Machine translation performance on WMT benchmarks. Scores are averaged across both translation directions. Standard deviations (approximately 1.0 \sim 2.0) are reported in Table 10 in Appendix C.

Model	Method	Metric	WMT14	WMT16	WMT16
			En-Fr	En-De	En-Ro
Phi-3.5-4B	Base	BLEU \uparrow	33.6	39.2	17.7
		CHRF \uparrow	58.3	63.2	45.5
		TER \downarrow	53.0	47.9	73.4
	Cont. finetune	BLEU \uparrow	36.5	42.3	16.7
		CHRF \uparrow	61.0	65.4	45.8
		TER \downarrow	51.5	44.9	79.7
Phi-3-14B	zip2zip	BLEU \uparrow	34.1	39.7	14.3
		CHRF \uparrow	59.4	64.5	44.2
		TER \downarrow	54.5	48.0	93.5
	Base	BLEU \uparrow	39.1	43.1	21.3
		CHRF \uparrow	62.6	65.6	51.0
		TER \downarrow	49.3	44.1	70.5
	Cont. finetune	BLEU \uparrow	38.9	48.4	21.8
		CHRF \uparrow	63.2	70.1	52.0
		TER \downarrow	48.8	39.8	68.3
	zip2zip	BLEU \uparrow	36.4	44.8	19.5
		CHRF \uparrow	62.8	68.1	50.1
		TER \downarrow	51.2	42.9	72.9

Efficient LZW Tokenization. zip2zip introduces an additional LZW compression step during inference and a decompression step at the end of generation. As a result, the efficiency of LZW-integrated tokenization is important to overall performance. To minimize overhead, we implemented a Rust-based zip2zip tokenizer that outperforms the Python version (see Figure 4) and matches the latency of HuggingFace’s fast BPE tokenizer.

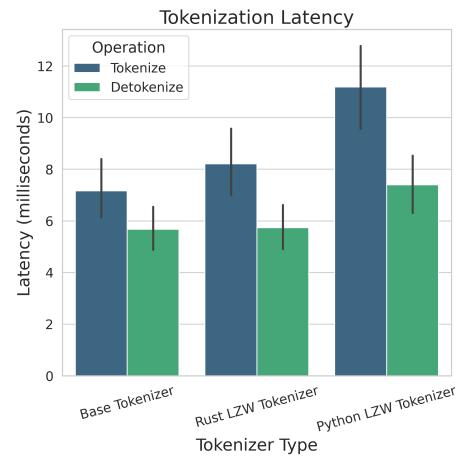


Figure 4: zip2zip tokenizer latency (ms) vs. HF tokenizer.

385
386
387
388
Table 6: **Throughput (tokens/sec)** comparison of the zip2zip framework against the baseline Hugging Face Transformers generate
389 and MLX generate implementation. Performance is detailed for prefill and decode phases across various context lengths (first value in
390 column headers) combined with a 256-token generation length. zip2zip demonstrates notable throughput improvements, in both prefill
391 and decoding phase.
392

Setting	Method	256+256		512+256		1024+256		2048+256	
		Prefill	Decode	Prefill	Decode	Prefill	Decode	Prefill	Decode
<i>Hardware: Apple M1 (16GB RAM)</i>									
Phi-3-4B	Base model	165.0	7.3	211.3	7.5	200.9	7.1	196.6	6.8
	zip2zip	145.5	7.9	231.4	10.1	189.6	7.4	233.8	7.3
	Relative %	-11.8%	+7.5%	+9.5%	+34.8%	-6.6%	+3.9%	+18.9%	+7.5%
<i>Hardware: NVIDIA H100 80GB GPU</i>									
Phi-3.5-4B	Base model	700.9	56.2	1347.2	54.4	2689.4	52.8	4993.2	53.1
	zip2zip	936.6	61.4	2722.1	79.8	4326.7	61.5	9258.1	61.9
	Relative %	+33.6%	+9.3%	+102.6%	+46.6%	+60.9%	+16.6%	+85.4%	+16.5%
Phi-3-14B	Base model	724.4	44.6	1356.3	43.8	2328.6	45.1	3849.5	42.2
	zip2zip	1024.6	54.9	1973.0	61.1	3657.0	66.8	7239.1	46.3
	Relative %	+41.5%	+23.0%	+45.5%	+39.5%	+57.0%	+48.1%	+88.1%	+9.6%

4. Related Work

Vocabulary Expansion. Several works have explored expanding the tokenizer vocabulary to better support specific domains or languages. Zhao et al. (2024); Kim et al. (2024); Liu et al. (2023; 2024) adapt LLaMA to Chinese, Korean, and specialized domains such as mental health and law by appending new tokens. Wang et al. (2025); Liu et al. (2024) conducted studies on how to effectively expand the vocabulary by better selecting the subset of tokens to add. In contrast, zip2zip is the first to enable *dynamic vocabulary expansion at inference time*, constructing new tokens based on the input context without requiring retraining or modifying the tokenizer ahead of time.

Prompt Compression. Prompt compression methods include GistTokens (Mu et al., 2023), Selective Context (Li et al., 2023), LLMLingua (Jiang et al., 2023), Summary Vectors (Chevalier et al., 2023), In-context Autoencoder (Ge et al., 2024), and others (Wingate et al., 2022) reduce the input token length and but do not impact the number of output tokens, which often dominates overall generation time. In contrast, zip2zip compresses *both* the input and output token sequences.

Latent Tokens Representation. The concept of latent token representations, or *patches*, has been mostly explored in computer vision, with methods like Token Merging (Bolya et al., 2023) and Token Pooling (Marin et al., 2023) aiming to reduce sequence length while preserving semantic content. Recently, Byte Latent Transformer (BLT) (Pagnoni et al., 2024) extended this concept to language modeling by discarding tokens entirely and operating directly at the byte level. Both BLT and zip2zip adopt a hierarchical modeling of input for LLMs, but they differ in three key

ways: (1) **Goal**: BLT aims to replace the tokenizer, whereas zip2zip seeks to expand and improve it; (2) **Algorithm**: BLT uses entropy-based segmentation, while zip2zip applies LZW-based token compression; (3) **Training**: BLT requires training from scratch, whereas zip2zip enables continued adaptation of pretrained models. Lester et al. (2024) propose improving language model efficiency by training LLMs directly on text compressed with arithmetic coding. While both approaches leverage compression to enhance efficiency, zip2zip emphasizes dynamic vocabulary expansion to enable uptraining of existing models. In contrast, Lester et al. (2024) requires training from scratch.

5. Discussion and Limitations

Beyond LZW. While we adopt LZW for dynamic construction of hypertokens, zip2zip is broadly compatible with any online compression algorithm. Future work may explore alternative schemes that provide different trade-offs between compression efficiency and model performance.

Codebook Management Strategy. The LZW algorithm grows the codebook linearly with the number of tokens in the context window. Empirical results show that only about 25% of hypertokens are reused during generation, leaving substantial room for optimization. Two potential improvements are: (1) **pruning** or **selective retention** strategies to reduce unused entries, and (2) **codebook prefilling**, which could be beneficial if likely tokens can be anticipated before input processing.

Compression-Quality Trade-off. There is an inherent trade-off between compression and modeling: as the token space is compressed more aggressively, redundancy is reduced—but so is predictability—making it harder for

the model to forecast the next (hyper)token. In the extreme, optimal compression schemes such as arithmetic coding produce sequences that are statistically indistinguishable from random noise, rendering them unlearnable by language models (Lester et al., 2024). Empirically, we observe this effect as increased perplexity under higher compression levels (Table 3), which can undermine the benefits of compression by degrading generation quality (though minor in the tasks in Table ?? and Table 5). Striking the right balance between compression and model performance remains an important direction for future research.

6. Conclusion

We introduced zip2zip, a framework for inference-time vocabulary adaptation with LLMs. By integrating LZW-based token compression with a dynamic hypertoken embedding mechanism, zip2zip enables substantial reductions in sequence length and decoding steps, leading to improved inference efficiency with minimal architectural modifications. Our experiments demonstrate that zip2zip maintains strong performance across a range of tasks while achieving significant gains in inference efficiency. These findings highlight the promise of integrating dynamic tokenization into LLMs, opening up new directions for research in LLM efficiency.

References

- Abdin, M., Aneja, J., Behl, H., Bubeck, S., Eldan, R., Guanasekar, S., Harrison, M., Hewett, R. J., Javaheripi, M., Kauffmann, P., Lee, J. R., Lee, Y. T., Li, Y., Liu, W., Mendes, C. C. T., Nguyen, A., Price, E., de Rosa, G., Saarikivi, O., Salim, A., Shah, S., Wang, X., Ward, R., Wu, Y., Yu, D., Zhang, C., and Zhang, Y. Phi-4 technical report, 2024. URL <https://arxiv.org/abs/2412.08905>.
- Ahia, O., Kumar, S., Gonen, H., Kasai, J., Mortensen, D., Smith, N., and Tsvetkov, Y. Do all languages cost the same? tokenization in the era of commercial language models. In Bouamor, H., Pino, J., and Bali, K. (eds.), *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pp. 9904–9923, Singapore, December 2023. Association for Computational Linguistics. doi: 10.18653/v1/2023.emnlp-main.614. URL <https://aclanthology.org/2023.emnlp-main.614/>.
- Bisk, Y., Zellers, R., Bras, R. L., Gao, J., and Choi, Y. Piqa: Reasoning about physical commonsense in natural language. In *AAAI Conference on Artificial Intelligence*, 2019. URL <https://api.semanticscholar.org/CorpusID:208290939>.
- Bojar, O., Chatterjee, R., Federmann, C., Graham, Y., Hadidow, B., Huck, M., Jimeno Yepes, A., Koehn, P., Logacheva, V., Monz, C., Negri, M., Névéol, A., Neves, M., Popel, M., Post, M., Rubino, R., Scarton, C., Specia, L., Turchi, M., Verspoor, K., and Zampieri, M. Findings of the 2016 conference on machine translation. In Bojar, O., Buck, C., Chatterjee, R., Federmann, C., Guillou, L., Haddow, B., Huck, M., Yepes, A. J., Névéol, A., Neves, M., Pecina, P., Popel, M., Koehn, P., Monz, C., Negri, M., Post, M., Specia, L., Verspoor, K., Tiedemann, J., and Turchi, M. (eds.), *Proceedings of the First Conference on Machine Translation: Volume 2, Shared Task Papers*, pp. 131–198, Berlin, Germany, August 2016. Association for Computational Linguistics. doi: 10.18653/v1/W16-2301. URL <https://aclanthology.org/W16-2301/>.
- Bolya, D., Fu, C.-Y., Dai, X., Zhang, P., Feichtenhofer, C., and Hoffman, J. Token merging: Your vit but faster, 2023. URL <https://arxiv.org/abs/2210.09461>.
- Brown, T., Mann, B., Ryder, N., Subbiah, M., Kaplan, J. D., Dhariwal, P., Neelakantan, A., Shyam, P., Sastry, G., Askell, A., Agarwal, S., Herbert-Voss, A., Krueger, G., Henighan, T., Child, R., Ramesh, A., Ziegler, D., Wu, J., Winter, C., Hesse, C., Chen, M., Sigler, E., Litwin, M., Gray, S., Chess, B., Clark, J., Berner, C., McCandlish, S., Radford, A., Sutskever, I., and Amodei, D. Language models are few-shot learners. In Larochelle, H., Ranzato, M., Hadsell, R., Balcan, M., and Lin, H. (eds.), *Advances in Neural Information Processing Systems*, volume 33, pp. 1877–1901. Curran Associates, Inc., 2020. URL https://proceedings.neurips.cc/paper_files/paper/2020/file/1457c0d6bfcb4967418fb8ac142f64a-Paper.pdf.
- Bubeck, S., Chandrasekaran, V., Eldan, R., Gehrke, J., Horvitz, E., Kamar, E., Lee, P., Lee, Y. T., Li, Y., Lundberg, S., Nori, H., Palangi, H., Ribeiro, M. T., and Zhang, Y. Sparks of artificial general intelligence: Early experiments with gpt-4, 2023. URL <https://arxiv.org/abs/2303.12712>.
- Chevalier, A., Wettig, A., Ajith, A., and Chen, D. Adapting language models to compress contexts. In Bouamor, H., Pino, J., and Bali, K. (eds.), *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pp. 3829–3846, Singapore, December 2023. Association for Computational Linguistics. doi: 10.18653/v1/2023.emnlp-main.232. URL <https://aclanthology.org/2023.emnlp-main.232/>.
- Clark, P., Cowhey, I., Etzioni, O., Khot, T., Sabharwal, A., Schoenick, C., and Tafjord, O. Think you have solved question answering? try arc, the ai2 reasoning challenge, 2018. URL <https://arxiv.org/abs/1803.05457>.

- 495 Cobbe, K., Kosaraju, V., Bavarian, M., Chen, M., Jun, H.,
496 Kaiser, L., Plappert, M., Tworek, J., Hilton, J., Nakano,
497 R., Hesse, C., and Schulman, J. Training verifiers to solve
498 math word problems. *arXiv preprint arXiv:2110.14168*,
499 2021.
- 500 Dagan, G., Synnaeve, G., and Rozière, B. Getting the
501 most out of your tokenizer for pre-training and domain
502 adaptation, 2024. URL <https://arxiv.org/abs/2402.01035>.
- 503 Ding, N., Chen, Y., Xu, B., Qin, Y., Zheng, Z., Hu, S., Liu,
504 Z., Sun, M., and Zhou, B. Enhancing chat language
505 models by scaling high-quality instructional conversations,
506 2023.
- 507 Frieder, S., Pinchetti, L., Chevalier, A., Griffiths, R.-R.,
508 Salvatori, T., Lukasiewicz, T., Petersen, P. C., and Berner,
509 J. Mathematical capabilities of chatgpt, 2023. URL
510 <https://arxiv.org/abs/2301.13867>.
- 511 Gao, L., Biderman, S., Black, S., Golding, L., Hoppe, T.,
512 Foster, C., Phang, J., He, H., Thite, A., Nabeshima, N.,
513 Presser, S., and Leahy, C. The pile: An 800gb dataset of
514 diverse text for language modeling, 2020. URL <https://arxiv.org/abs/2101.00027>.
- 515 Gao, L., Tow, J., Abbasi, B., Biderman, S., Black, S., DiPofi,
516 A., Foster, C., Golding, L., Hsu, J., Le Noac'h, A., Li,
517 H., McDonell, K., Muennighoff, N., Ociepa, C., Phang,
518 J., Reynolds, L., Schoelkopf, H., Skowron, A., Sutawika,
519 L., Tang, E., Thite, A., Wang, B., Wang, K., and Zou, A.
520 The language model evaluation harness, 07 2024. URL
521 <https://zenodo.org/records/12608602>.
- 522 Ge, T., Jing, H., Wang, L., Wang, X., Chen, S.-Q., and
523 Wei, F. In-context autoencoder for context compression
524 in a large language model. In *The Twelfth International
525 Conference on Learning Representations*, 2024. URL
526 <https://openreview.net/forum?id=uREj4ZuGJE>.
- 527 Grattafiori, A., Dubey, A., Jauhri, A., Pandey, A., Kadian,
528 A., Al-Dahle, A., Letman, A., Mathur, A., Schelten, A.,
529 Vaughan, A., Yang, A., Fan, A., Goyal, A., Hartshorn,
530 A., Yang, A., Mitra, A., Srivankumar, A., Korenev,
531 A., Hinsvark, A., Rao, A., Zhang, A., Rodriguez, A.,
532 Gregerson, A., Spataru, A., Roziere, B., Biron, B., Tang,
533 B., Chern, B., Caucheteux, C., Nayak, C., Bi, C., Marra,
534 C., McConnell, C., Keller, C., Touret, C., Wu, C., Wong,
535 C., Ferrer, C. C., Nikolaidis, C., Allonsius, D., Song, D.,
536 Pintz, D., Livshits, D., Wyatt, D., Esiobu, D., Choudhary,
537 D., Mahajan, D., Garcia-Olano, D., Perino, D., Hupkes,
538 D., Lakomkin, E., AlBadawy, E., Lobanova, E., Dinan,
539 E., Smith, E. M., Radenovic, F., Guzmán, F., Zhang, F.,
540 Synnaeve, G., Lee, G., Anderson, G. L., Thattai, G., Nail,
541 G., Mialon, G., Pang, G., Cucurell, G., Nguyen, H., Ko-
542 revar, H., Xu, H., Touvron, H., Zarov, I., Ibarra, I. A.,
543 Kloumann, I., Misra, I., Evtimov, I., Zhang, J., Copet, J.,
544 Lee, J., Geffert, J., Vranes, J., Park, J., Mahadeokar, J.,
545 Shah, J., van der Linde, J., Billock, J., Hong, J., Lee, J.,
546 Fu, J., Chi, J., Huang, J., Liu, J., Wang, J., Yu, J., Bitton,
547 J., Spisak, J., Park, J., Rocca, J., Johnstun, J., Saxe, J., Jia,
548 J., Alwala, K. V., Prasad, K., Upasani, K., Plawiak, K., Li,
549 K., Heafield, K., Stone, K., El-Arini, K., Iyer, K., Malik,
550 K., Chiu, K., Bhalla, K., Lakhota, K., Rantala-Yeary,
551 L., van der Maaten, L., Chen, L., Tan, L., Jenkins, L.,
552 Martin, L., Madaan, L., Malo, L., Blecher, L., Landzaat,
553 L., de Oliveira, L., Muzzi, M., Pasupuleti, M., Singh,
554 M., Paluri, M., Kardas, M., Tsimpoukelli, M., Oldham,
555 M., Rita, M., Pavlova, M., Kambadur, M., Lewis, M.,
556 Si, M., Singh, M. K., Hassan, M., Goyal, N., Torabi, N.,
557 Bashlykov, N., Bogoychev, N., Chatterji, N., Zhang, N.,
558 Duchenne, O., Celebi, O., Alrassy, P., Zhang, P., Li, P.,
559 Vasic, P., Weng, P., Bhargava, P., Dubal, P., Krishnan,
560 P., Koura, P. S., Xu, P., He, Q., Dong, Q., Srinivasan,
561 R., Ganapathy, R., Calderer, R., Cabral, R. S., Stojnic,
562 R., Raileanu, R., Maheswari, R., Girdhar, R., Patel, R.,
563 Sauvestre, R., Polidoro, R., Sumbaly, R., Taylor, R., Silva,
564 R., Hou, R., Wang, R., Hosseini, S., Chennabasappa, S.,
565 Singh, S., Bell, S., Kim, S. S., Edunov, S., Nie, S., Narang,
566 S., Raparthy, S., Shen, S., Wan, S., Bhosale, S., Zhang,
567 S., Vandenhende, S., Batra, S., Whitman, S., Sootla, S.,
568 Collot, S., Gururangan, S., Borodinsky, S., Herman, T.,
569 Fowler, T., Sheasha, T., Georgiou, T., Scialom, T., Speck-
570 bacher, T., Mihaylov, T., Xiao, T., Karn, U., Goswami, V.,
571 Gupta, V., Ramanathan, V., Kerkez, V., Gonguet, V., Do,
572 V., Vogeti, V., Albiero, V., Petrovic, V., Chu, W., Xiong,
573 W., Fu, W., Meers, W., Martinet, X., Wang, X., Wang,
574 X., Tan, X. E., Xia, X., Xie, X., Jia, X., Wang, X., Gold-
575 schlag, Y., Gaur, Y., Babaei, Y., Wen, Y., Song, Y., Zhang,
576 Y., Li, Y., Mao, Y., Couder, Z. D., Yan, Z., Chen, Z.,
577 Papakipos, Z., Singh, A., Srivastava, A., Jain, A., Kelsey,
578 A., Shajnfeld, A., Gangidi, A., Victoria, A., Goldstand,
579 A., Menon, A., Sharma, A., Boesenberg, A., Baevski, A.,
580 Feinstein, A., Kallet, A., Sangani, A., Teo, A., Yunus, A.,
581 Lupu, A., Alvarado, A., Caples, A., Gu, A., Ho, A., Poult-
582 ton, A., Ryan, A., Ramchandani, A., Dong, A., Franco,
583 A., Goyal, A., Saraf, A., Chowdhury, A., Gabriel, A.,
584 Bharambe, A., Eisenman, A., Yazdan, A., James, B.,
585 Maurer, B., Leonhardi, B., Huang, B., Loyd, B., Paola,
586 B. D., Paranjape, B., Liu, B., Wu, B., Ni, B., Hancock,
587 B., Wasti, B., Spence, B., Stojkovic, B., Gamido, B.,
588 Montalvo, B., Parker, C., Burton, C., Mejia, C., Liu, C.,
589 Wang, C., Kim, C., Zhou, C., Hu, C., Chu, C.-H., Cai, C.,
590 Tindal, C., Feichtenhofer, C., Gao, C., Civin, D., Beaty,
591 D., Kreymer, D., Li, D., Adkins, D., Xu, D., Testuggine,
592 D., David, D., Parikh, D., Liskovich, D., Foss, D., Wang,
593 D., Le, D., Holland, D., Dowling, E., Jamil, E., Mont-
594 gomery, E., Presani, E., Hahn, E., Wood, E., Le, E.-T.,
595 Brinkman, E., Arcaute, E., Dunbar, E., Smothers, E., Sun,
596 F., Kreuk, F., Tian, F., Kokkinos, F., Ozgenel, F., Cag-

550 gioni, F., Kanayet, F., Seide, F., Florez, G. M., Schwarz,
551 G., Badeer, G., Swee, G., Halpern, G., Herman, G., Sizov,
552 G., Guangyi, Zhang, Lakshminarayanan, G., Inan, H.,
553 Shojanazeri, H., Zou, H., Wang, H., Zha, H., Habeeb, H.,
554 Rudolph, H., Suk, H., Aspegren, H., Goldman, H., Zhan,
555 H., Damlaj, I., Molybog, I., Tufanov, I., Leontiadis, I.,
556 Veliche, I.-E., Gat, I., Weissman, J., Geboski, J., Kohli,
557 J., Lam, J., Asher, J., Gaya, J.-B., Marcus, J., Tang, J.,
558 Chan, J., Zhen, J., Reizenstein, J., Teboul, J., Zhong, J.,
559 Jin, J., Yang, J., Cummings, J., Carvill, J., Shepard, J.,
560 McPhie, J., Torres, J., Ginsburg, J., Wang, J., Wu, K., U,
561 K. H., Saxena, K., Khandelwal, K., Zand, K., Matosich,
562 K., Veeraraghavan, K., Michelena, K., Li, K., Jagadeesh,
563 K., Huang, K., Chawla, K., Huang, K., Chen, L., Garg,
564 L., A, L., Silva, L., Bell, L., Zhang, L., Guo, L., Yu, L.,
565 Moshkovich, L., Wehrstedt, L., Khabsa, M., Avalani, M.,
566 Bhatt, M., Mankus, M., Hasson, M., Lennie, M., Reso,
567 M., Groshev, M., Naumov, M., Lathi, M., Keneally, M.,
568 Liu, M., Seltzer, M. L., Valko, M., Restrepo, M., Patel,
569 M., Vyatskov, M., Samvelyan, M., Clark, M., Macey,
570 M., Wang, M., Hermoso, M. J., Metanat, M., Rastegari,
571 M., Bansal, M., Santhanam, N., Parks, N., White, N.,
572 Bawa, N., Singhal, N., Egebo, N., Usunier, N., Mehta,
573 N., Laptev, N. P., Dong, N., Cheng, N., Chernoguz, O.,
574 Hart, O., Salpekar, O., Kalinli, O., Kent, P., Parekh, P.,
575 Saab, P., Balaji, P., Rittner, P., Bontrager, P., Roux, P.,
576 Dollar, P., Zvyagina, P., Ratanchandani, P., Yuvraj, P.,
577 Liang, Q., Alao, R., Rodriguez, R., Ayub, R., Murthy, R.,
578 Nayani, R., Mitra, R., Parthasarathy, R., Li, R., Hogan,
579 R., Battey, R., Wang, R., Howes, R., Rinott, R., Mehta,
580 S., Siby, S., Bondu, S. J., Datta, S., Chugh, S., Hunt, S.,
581 Dhillon, S., Sidorov, S., Pan, S., Mahajan, S., Verma,
582 S., Yamamoto, S., Ramaswamy, S., Lindsay, S., Lindsay,
583 S., Feng, S., Lin, S., Zha, S. C., Patil, S., Shankar, S.,
584 Zhang, S., Zhang, S., Wang, S., Agarwal, S., Sajuyigbe,
585 S., Chintala, S., Max, S., Chen, S., Kehoe, S., Satterfield,
586 S., Govindaprasad, S., Gupta, S., Deng, S., Cho,
587 S., Virk, S., Subramanian, S., Choudhury, S., Goldman,
588 S., Remez, T., Glaser, T., Best, T., Koehler, T., Robinson,
589 T., Li, T., Zhang, T., Matthews, T., Chou, T., Shaked,
590 T., Vontimitta, V., Ajayi, V., Montanez, V., Mohan, V.,
591 Kumar, V. S., Mangla, V., Ionescu, V., Poenaru, V., Mi-
592 hailescu, V. T., Ivanov, V., Li, W., Wang, W., Jiang, W.,
593 Bouaziz, W., Constable, W., Tang, X., Wu, X., Wang, X.,
594 Wu, X., Gao, X., Kleinman, Y., Chen, Y., Hu, Y., Jia, Y.,
595 Qi, Y., Li, Y., Zhang, Y., Zhang, Y., Adi, Y., Nam, Y., Yu,
596 Wang, Zhao, Y., Hao, Y., Qian, Y., Li, Y., He, Y., Rait,
597 Z., DeVito, Z., Rosnbrick, Z., Wen, Z., Yang, Z., Zhao,
598 Z., and Ma, Z. The llama 3 herd of models, 2024. URL
599 <https://arxiv.org/abs/2407.21783>.

600
601 Hu, E. J., Shen, Y., Wallis, P., Allen-Zhu, Z., Li, Y.,
602 Wang, S., Wang, L., and Chen, W. LoRA: Low-rank
603 adaptation of large language models. In *International*

604 Conference on Learning Representations, 2022. URL
605 <https://openreview.net/forum?id=nZeVKeeFYf9>.

606 Jiang, H., Wu, Q., Lin, C.-Y., Yang, Y., and Qiu, L. LLM-
607 Lingua: Compressing prompts for accelerated inference
608 of large language models. In Bouamor, H., Pino, J., and
609 Bali, K. (eds.), *Proceedings of the 2023 Conference on
610 Empirical Methods in Natural Language Processing*, pp.
611 13358–13376, Singapore, December 2023. Association
612 for Computational Linguistics. doi: 10.18653/v1/2023.
613 emnlp-main.825. URL <https://aclanthology.org/2023.emnlp-main.825>.

614 Jiang, J., Wang, F., Shen, J., Kim, S., and Kim, S. A survey
615 on large language models for code generation, 2024. URL
616 <https://arxiv.org/abs/2406.00515>.

617 Kaplan, J., McCandlish, S., Henighan, T., Brown, T. B.,
618 Chess, B., Child, R., Gray, S., Radford, A., Wu, J., and
619 Amodei, D. Scaling laws for neural language models,
620 2020. URL <https://arxiv.org/abs/2001.08361>.

621 Kim, S., Choi, S., and Jeong, M. Efficient and effective vo-
622 cabulary expansion towards multilingual large language
623 models, 2024. URL <https://arxiv.org/abs/2402.14714>.

624 Lester, B., Lee, J., Alemi, A., Pennington, J., Roberts, A.,
625 Sohl-Dickstein, J., and Constant, N. Training llms over
626 neurally compressed text, 2024. URL <https://arxiv.org/abs/2404.03626>.

627 LI, J., Beeching, E., Tunstall, L., Lipkin, B., Soletskyi,
628 R., Huang, S. C., Rasul, K., Yu, L., Jiang, A., Shen, Z.,
629 Qin, Z., Dong, B., Zhou, L., Fleureau, Y., Lample, G.,
630 and Polu, S. Numinamath. [<https://huggingface.co/AI-MO/NuminaMath-1.5>] (https://github.com/project-numina/aimo-progress-prize/blob/main/report/numina_dataset.pdf), 2024.

631 Li, Y., Dong, B., Guerin, F., and Lin, C. Compressing
632 context to enhance inference efficiency of large language
633 models. In Bouamor, H., Pino, J., and Bali, K. (eds.), *Pro-
634 ceedings of the 2023 Conference on Empirical Methods
635 in Natural Language Processing*, pp. 6342–6353, Singa-
636 pore, December 2023. Association for Computational
637 Linguistics. doi: 10.18653/v1/2023.emnlp-main.391.
638 URL <https://aclanthology.org/2023.emnlp-main.391>.

639 Liang, D., Gonen, H., Mao, Y., Hou, R., Goyal, N.,
640 Ghazvininejad, M., Zettlemoyer, L., and Khabsa, M. XLM-V:
641 Overcoming the vocabulary bottleneck in multi-
642 lingual masked language models. In Bouamor, H., Pino,
643 J., and Bali, K. (eds.), *Proceedings of the 2023 Confer-
644 ence on Empirical Methods in Natural Language Process-
645 ing*, pp. 13142–13152, Singapore, December 2023. Asso-
646 ciation for Computational Linguistics. doi: 10.18653/v1/

- 605 2023.emnlp-main.813. URL <https://aclanthology.org/2023.emnlp-main.813/>.
- 606
- 607 Liu, C., Wang, S., Qing, L., Kuang, K., Kang, Y., Sun,
608 C., and Wu, F. Gold panning in vocabulary: An
609 adaptive method for vocabulary expansion of domain-
610 specific LLMs. In Al-Onaizan, Y., Bansal, M., and
611 Chen, Y.-N. (eds.), *Proceedings of the 2024 Confer-
612 ence on Empirical Methods in Natural Language Pro-
613 cessing*, pp. 7442–7459, Miami, Florida, USA, No-
614 vember 2024. Association for Computational Linguis-
615 tics. doi: 10.18653/v1/2024.emnlp-main.424. URL <https://aclanthology.org/2024.emnlp-main.424/>.
- 616
- 617 Liu, S., Deng, N., Sabour, S., Jia, Y., Huang, M., and Miha-
618 lea, R. Task-adaptive tokenization: Enhancing long-form
619 text generation efficacy in mental health and beyond. In
620 Bouamor, H., Pino, J., and Bali, K. (eds.), *Proceedings
621 of the 2023 Conference on Empirical Methods in Natural
622 Language Processing*, pp. 15264–15281, Singapore, De-
623 cember 2023. Association for Computational Linguis-
624 tics. doi: 10.18653/v1/2023.emnlp-main.944. URL <https://aclanthology.org/2023.emnlp-main.944/>.
- 625
- 626 Lozhkov, A., Ben Allal, L., von Werra, L., and Wolf, T.
627 Fineweb-edu: the finest collection of educational content,
628 2024a. URL <https://huggingface.co/datasets/HuggingFaceFW/fineweb-edu>.
- 629
- 630 Lozhkov, A., Li, R., Allal, L. B., Cassano, F., Lamy-Poirier,
631 J., Tazi, N., Tang, A., Pykhtar, D., Liu, J., Wei, Y., Liu, T.,
632 Tian, M., Kocetkov, D., Zucker, A., Belkada, Y., Wang,
633 Z., Liu, Q., Abulkhanov, D., Paul, I., Li, Z., Li, W.-D.,
634 Risdal, M., Li, J., Zhu, J., Zhuo, T. Y., Zheltonozhskii,
635 E., Dade, N. O. O., Yu, W., Krauß, L., Jain, N., Su, Y.,
636 He, X., Dey, M., Abati, E., Chai, Y., Muennighoff, N.,
637 Tang, X., Oblokulov, M., Akiki, C., Marone, M., Mou,
638 C., Mishra, M., Gu, A., Hui, B., Dao, T., Zebaze, A.,
639 Dehaene, O., Patry, N., Xu, C., McAuley, J., Hu, H.,
640 Scholak, T., Paquet, S., Robinson, J., Anderson, C. J.,
641 Chapados, N., Patwary, M., Tajbakhsh, N., Jernite, Y.,
642 Ferrandis, C. M., Zhang, L., Hughes, S., Wolf, T., Guha,
643 A., von Werra, L., and de Vries, H. Starcoder 2 and the
644 stack v2: The next generation, 2024b.
- 645
- 646 Macháček, M. and Bojar, O. Results of the WMT14 met-
647 rics shared task. In Bojar, O., Buck, C., Federmann, C.,
648 Haddow, B., Koehn, P., Monz, C., Post, M., and Spe-
649 cia, L. (eds.), *Proceedings of the Ninth Workshop on
650 Statistical Machine Translation*, pp. 293–301, Baltimore,
651 Maryland, USA, June 2014. Association for Compu-
652 tational Linguistics. doi: 10.3115/v1/W14-3336. URL
653 <https://aclanthology.org/W14-3336/>.
- 654
- 655 Magnusson, I., Bhagia, A., Hofmann, V., Soldaini, L.,
656 Jha, A., Tafjord, O., Schwenk, D., Walsh, P., Elazar,
657
- 658 Y., Lo, K., Groeneveld, D., Beltagy, I., Hajishirzi,
659 H., Smith, N. A., Richardson, K., and Dodge, J.
660 Paloma: A benchmark for evaluating language model
661 fit. *ArXiv*, abs/2312.10523, 2023. URL <https://api.semanticscholar.org/CorpusID:266348815>.
- 662
- 663 Marin, D., Chang, J.-H. R., Ranjan, A., Prabhu, A., Raste-
664 gari, M., and Tuzel, O. Token pooling in vision trans-
665 formers for image classification. In *2023 IEEE/CVF Winter
666 Conference on Applications of Computer Vision (WACV)*,
667 pp. 12–21, 2023. doi: 10.1109/WACV56688.2023.00010.
- 668
- 669 Merity, S., Xiong, C., Bradbury, J., and Socher, R. Pointer
670 sentinel mixture models, 2016.
- 671
- 672 Mihaylov, T., Clark, P., Khot, T., and Sabharwal, A. Can
673 a suit of armor conduct electricity? a new dataset for
674 open book question answering. In Riloff, E., Chiang,
675 D., Hockenmaier, J., and Tsujii, J. (eds.), *Proceedings
676 of the 2018 Conference on Empirical Methods in Natural
677 Language Processing*, pp. 2381–2391, Brussels, Belgium,
678 October–November 2018. Association for Computational
679 Linguistics. doi: 10.18653/v1/D18-1260. URL <https://aclanthology.org/D18-1260/>.
- 680
- 681 Mu, J., Li, X., and Goodman, N. Learning to compress
682 prompts with gist tokens. In Oh, A., Naumann, T.,
683 Globerson, A., Saenko, K., Hardt, M., and Levine, S.
684 (eds.), *Advances in Neural Information Processing
685 Systems*, volume 36, pp. 19327–19352. Curran As-
686 sociates, Inc., 2023. URL https://proceedings.neurips.cc/paper_files/paper/2023/file/3d77c6dcc7f143aa2154e7f4d5e22d68-Paper-Conference.pdf.
- 687
- 688 Nori, H., King, N., McKinney, S. M., Carignan, D., and
689 Horvitz, E. Capabilities of gpt-4 on medical challenge
690 problems, 2023. URL <https://arxiv.org/abs/2303.13375>.
- 691
- 692 Pagnoni, A., Pasunuru, R., Rodriguez, P., Nguyen, J.,
693 Muller, B., Li, M., Zhou, C., Yu, L., Weston, J., Zettle-
694 moyer, L., Ghosh, G., Lewis, M., Holtzman, A., and
695 Iyer, S. Byte latent transformer: Patches scale better
696 than tokens, 2024. URL <https://arxiv.org/abs/2412.09871>.
- 697
- 698 Paperno, D., Kruszewski, G., Lazaridou, A., Pham, N. Q.,
699 Bernardi, R., Pezzelle, S., Baroni, M., Boleda, G., and
700 Fernández, R. The LAMBADA dataset: Word prediction
701 requiring a broad discourse context. In Erk, K. and Smith,
702 N. A. (eds.), *Proceedings of the 54th Annual Meeting of
703 the Association for Computational Linguistics (Volume 1:
704 Long Papers)*, pp. 1525–1534, Berlin, Germany, August
705 2016. Association for Computational Linguistics. doi:
706 10.18653/v1/P16-1144. URL [https://aclanthology.org/P16-1144/](https://aclanthology.org/P16-1144).

- 660 Penedo, G., Kydlíček, H., Sabolčec, V., Messmer, B.,
661 Foroutan, N., Jaggi, M., von Werra, L., and Wolf,
662 T. Fineweb2: A sparkling update with 1000s of
663 languages, 2024. URL <https://huggingface.co/datasets/HuggingFaceFW/fineweb-2>.
- 664
- 665 Petrov, A., La Malfa, E., H. S. Torr, P., and Bibi, A. Language model tokenizers introduce unfairness between languages. In *Advances in Neural Information Processing Systems*, 2023. URL <https://arxiv.org/abs/2305.15425>.
- 666
- 667 Radford, A., Wu, J., Child, R., Luan, D., Amodei, D.,
668 and Sutskever, I. Language models are unsupervised
669 multitask learners. *OpenAI blog*, 2019. URL <https://api.semanticscholar.org/CorpusID:160025533>.
- 670
- 671 Sakaguchi, K., Bras, R. L., Bhagavatula, C., and Choi,
672 Y. Winogrande. *Communications of the ACM*, 64:99 –
673 106, 2019. URL <https://api.semanticscholar.org/CorpusID:198893658>.
- 674
- 675 Sennrich, R., Haddow, B., and Birch, A. Neural machine
676 translation of rare words with subword units. In Erk,
677 K. and Smith, N. A. (eds.), *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 1715–1725,
678 Berlin, Germany, August 2016. Association for Computational Linguistics. doi: 10.18653/v1/P16-1162. URL
679 <https://aclanthology.org/P16-1162/>.
- 680
- 681 Singh, A. K. and Strouse, D. Tokenization counts: the
682 impact of tokenization on arithmetic in frontier llms, 2024.
683 URL <https://arxiv.org/abs/2402.14903>.
- 684
- 685 Sutskever, I., Vinyals, O., and Le, Q. V. Sequence to se-
686 quence learning with neural networks. In *Proceedings of the 28th International Conference on Neural Information Processing Systems - Volume 2*, NIPS’14, pp. 3104–3112,
687 Cambridge, MA, USA, 2014. MIT Press.
- 688
- 689 Team, G., Kamath, A., Ferret, J., Pathak, S., Vieillard,
690 N., Merhej, R., Perrin, S., Matejovicova, T., Ramé, A.,
691 Rivière, M., Rouillard, L., Mesnard, T., Cideron, G.,
692 bastien Grill, J., Ramos, S., Yvinec, E., Casbon, M., Pot,
693 E., Penchev, I., Liu, G., Visin, F., Kenealy, K., Beyer,
694 L., Zhai, X., Tsitsulin, A., Busa-Fekete, R., Feng, A.,
695 Sachdeva, N., Coleman, B., Gao, Y., Mustafa, B., Barr, I.,
696 Parisotto, E., Tian, D., Eyal, M., Cherry, C., Peter, J.-T.,
697 Sinopalnikov, D., Bhupatiraju, S., Agarwal, R., Kazemi,
698 M., Malkin, D., Kumar, R., Vilar, D., Brusilovsky, I.,
699 Luo, J., Steiner, A., Friesen, A., Sharma, A., Sharma,
700 A., Gilady, A. M., Goedekemeyer, A., Saade, A., Feng,
701 A., Kolesnikov, A., Bendebury, A., Abdagic, A., Vadi,
702 A., György, A., Pinto, A. S., Das, A., Bapna, A., Miech,
703 A., Yang, A., Paterson, A., Shenoy, A., Chakrabarti, A.,
704 Piot, B., Wu, B., Shahriari, B., Petrini, B., Chen, C.,
705 Lan, C. L., Choquette-Choo, C. A., Carey, C., Brick, C.,
706 Deutsch, D., Eisenbud, D., Cattle, D., Cheng, D., Pa-
707 paras, D., Sreepathihalli, D. S., Reid, D., Tran, D., Zelle,
708 D., Noland, E., Huizenga, E., Kharitonov, E., Liu, F.,
709 Amirkhanyan, G., Cameron, G., Hashemi, H., Klimczak-
710 Plucińska, H., Singh, H., Mehta, H., Lehri, H. T., Haz-
711 imeh, H., Ballantyne, I., Szpektor, I., Nardini, I., Pouget-
712 Abadie, J., Chan, J., Stanton, J., Wieting, J., Lai, J., Or-
713 bay, J., Fernandez, J., Newlan, J., yeong Ji, J., Singh,
714 J., Black, K., Yu, K., Hui, K., Vodrahalli, K., Greff, K.,
715 Qiu, L., Valentine, M., Coelho, M., Ritter, M., Hoffman,
716 M., Watson, M., Chaturvedi, M., Moynihan, M., Ma,
717 M., Babar, N., Noy, N., Byrd, N., Roy, N., Momchev,
718 N., Chauhan, N., Sachdeva, N., Bunyan, O., Botarda,
719 P., Caron, P., Rubenstein, P. K., Culliton, P., Schmid, P.,
720 Sessa, P. G., Xu, P., Stanczyk, P., Tafti, P., Shivanna, R.,
721 Wu, R., Pan, R., Rokni, R., Willoughby, R., Vallu, R.,
722 Mullins, R., Jerome, S., Smoot, S., Girgin, S., Iqbal, S.,
723 Reddy, S., Sheth, S., Pöder, S., Bhatnagar, S., Panyam,
724 S. R., Eiger, S., Zhang, S., Liu, T., Yacovone, T., Liechty,
725 T., Kalra, U., Evci, U., Misra, V., Roseberry, V., Feinberg,
726 V., Kolesnikov, V., Han, W., Kwon, W., Chen, X., Chow,
727 Y., Zhu, Y., Wei, Z., Egyed, Z., Cotrata, V., Giang, M.,
728 Kirk, P., Rao, A., Black, K., Babar, N., Lo, J., Moreira, E.,
729 Martins, L. G., Sanseviero, O., Gonzalez, L., Gleicher, Z.,
730 Warkentin, T., Mirrokni, V., Senter, E., Collins, E., Barral,
731 J., Ghahramani, Z., Hadsell, R., Matias, Y., Sculley, D.,
732 Petrov, S., Fiedel, N., Shazeer, N., Vinyals, O., Dean, J.,
733 Hassabis, D., Kavukcuoglu, K., Farabet, C., Buchatskaya,
734 E., Alayrac, J.-B., Anil, R., Dmitry, Lepikhin, Borgeaud,
735 S., Bachem, O., Joulin, A., Andreev, A., Hardin, C.,
736 Dadashi, R., and Huszenot, L. Gemma 3 technical report,
737 2025. URL <https://arxiv.org/abs/2503.19786>.
- 738
- Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J.,
739 Jones, L., Gomez, A. N., Kaiser, L. u., and Polosukhin,
740 I. Attention is all you need. In Guyon, I., Luxburg,
741 U. V., Bengio, S., Wallach, H., Fergus, R., Vish-
742 wanathan, S., and Garnett, R. (eds.), *Advances in Neural Information Processing Systems*, volume 30. Curran As-
743 sociates, Inc., 2017. URL https://proceedings.neurips.cc/paper_files/paper/2017/file/3f5ee243547dee91fdb053c1c4a845aa-Paper.pdf.
- 744
- Wang, H., Yu, D., Sun, K., Chen, J., and Yu, D. Improving
745 pre-trained multilingual model with vocabulary expansion.
746 In Bansal, M. and Villavicencio, A. (eds.), *Pro-
747 ceedings of the 23rd Conference on Computational Na-
748 tural Language Learning (CoNLL)*, pp. 316–327, Hong
749 Kong, China, November 2019. Association for Compu-
750 tational Linguistics. doi: 10.18653/v1/K19-1030. URL
751 <https://aclanthology.org/K19-1030/>.
- 752
- 753 Wang, S., Xie, Y., Ding, B., Gao, J., and Zhang, Y. Lan-
754 guage adaptation of large language models: An em-

715 pirical study on LLaMA2. In Rambow, O., Wanner,
716 L., Apidianaki, M., Al-Khalifa, H., Eugenio, B. D.,
717 and Schockaert, S. (eds.), *Proceedings of the 31st In-*
718 *ternational Conference on Computational Linguistics*,
719 pp. 7195–7208, Abu Dhabi, UAE, January 2025. As-
720 sociation for Computational Linguistics. URL <https://aclanthology.org/2025.coling-main.480/>.

721
722 Welch. A technique for high-performance data compression.
723 *Computer*, 17(6):8–19, 1984. doi: 10.1109/MC.1984.
724 1659158.

725
726 Wingate, D., Shoeybi, M., and Sorensen, T. Prompt com-
727 pression and contrastive conditioning for controllability
728 and toxicity reduction in language models. In Goldberg,
729 Y., Kozareva, Z., and Zhang, Y. (eds.), *Findings of the As-*
730 *sociation for Computational Linguistics: EMNLP 2022*,
731 pp. 5621–5634, Abu Dhabi, United Arab Emirates, De-
732 cember 2022. Association for Computational Linguistics.
733 doi: 10.18653/v1/2022.findings-emnlp.412. URL <https://aclanthology.org/2022.findings-emnlp.412/>.

734
735 Yang, A., Yang, B., Hui, B., Zheng, B., Yu, B., Zhou, C.,
736 Li, C., Li, C., Liu, D., Huang, F., Dong, G., Wei, H., Lin,
737 H., Tang, J., Wang, J., Yang, J., Tu, J., Zhang, J., Ma, J.,
738 Xu, J., Zhou, J., Bai, J., He, J., Lin, J., Dang, K., Lu, K.,
739 Chen, K., Yang, K., Li, M., Xue, M., Ni, N., Zhang, P.,
740 Wang, P., Peng, R., Men, R., Gao, R., Lin, R., Wang, S.,
741 Bai, S., Tan, S., Zhu, T., Li, T., Liu, T., Ge, W., Deng,
742 X., Zhou, X., Ren, X., Zhang, X., Wei, X., Ren, X., Fan,
743 Y., Yao, Y., Zhang, Y., Wan, Y., Chu, Y., Liu, Y., Cui, Z.,
744 Zhang, Z., and Fan, Z. Qwen2 technical report. *arXiv*
745 preprint arXiv:2407.10671, 2024.

746
747 Zellers, R., Holtzman, A., Bisk, Y., Farhadi, A., and Choi,
748 Y. HellaSwag: Can a machine really finish your sen-
749 tence? In Korhonen, A., Traum, D., and Màrquez,
750 L. (eds.), *Proceedings of the 57th Annual Meeting of*
751 *the Association for Computational Linguistics*, pp. 4791–
752 4800, Florence, Italy, July 2019. Association for Compu-
753 tational Linguistics. doi: 10.18653/v1/P19-1472. URL
754 <https://aclanthology.org/P19-1472/>.

755
756 Zhao, J., Zhang, Z., Gao, L., Zhang, Q., Gui, T., and
757 Huang, X. Llama beyond english: An empirical study
758 on language capability transfer, 2024. URL <https://arxiv.org/abs/2401.01055>.

770 A. Ablation Studies

771
772 **Definition A.1** (Compression Rate). We define the compression rate as the ratio between the number of tokens after
773 compression (N_{comp}) and the number of tokens in the original uncompressed text (N_{orig}), expressed as a percentage:

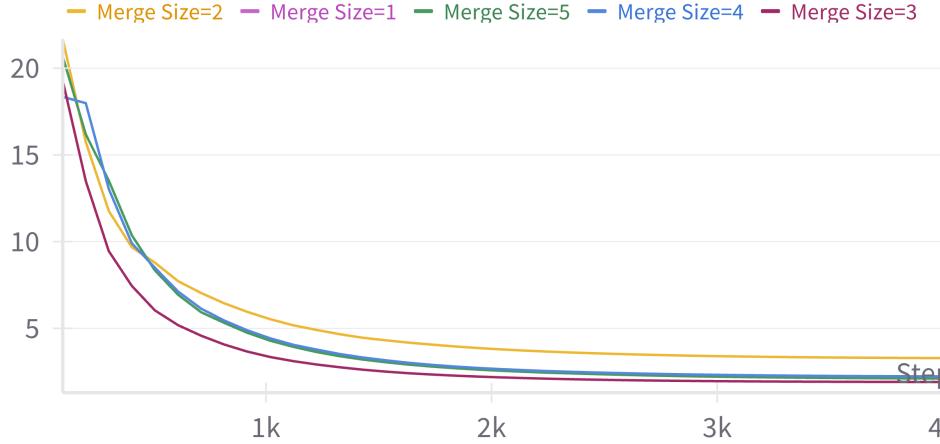
$$774 \quad \text{Compression Rate} = \frac{N_{comp}}{N_{orig}} \times 100\%.$$

775
776
777 A lower compression rate indicates greater reduction in token count, and thus more effective compression.

780 A.1. LZW Maximum Merge Size

781 The last column of Table 7 shows how the maximum merge size M affects compression rate when the context window
782 length is 2048. As M increases, compression rate improves significantly, especially from $M = 1$ to $M = 3$. Beyond that,
783 gains diminish, suggesting $M = 3$ strikes a good balance between efficiency and compression rate.

784
785 Interestingly, the relationship between maximum merge
786 size and training loss in Figure 5 as well as perplexity
787 in Table 7 is non-monotonic. The baseline case with
788 $M = 1$ (i.e., no zip2zip compression) yields the
789 lowest perplexity overall, which is expected and consistent
790 with prior findings that compression typically incurs a
791 trade-off in model performance. Among the compressed
792 settings, the case $M = 2$ performs the worst, with notice-
793 ably slower convergence and higher final loss. In contrast,
794 the case $M = 3$ achieves the best performance within
795 the compressed configurations, striking a favorable
796 balance between compression and prediction performance.
797 While $M = 4$ and $M = 5$ also perform reasonably well,
798 they exhibit slightly higher loss than $M = 3$, suggesting
799 diminishing returns or possible over-compression at larger maximum merge sizes (see Figure 5).



800
801
802 Figure 5: Effect of maximum merge size M on zip2zip training loss: $M = 1$ (no compression) achieves the lowest loss overall.
803 Among compressed settings, $M = 3$ performs best, while $M = 2$ shows the worst convergence. Larger M (4 and 5) yield slightly worse
804 results than $M = 3$.

805
806
807
808
809
810
811
812
813
814 Table 7 reports the byte-level perplexity across four corpora using a 1024-token context window. The results align closely
815 with the training loss trends observed earlier. Setting $M = 1$ (i.e., no compression) consistently achieves the lowest
816 perplexity across all datasets, reaffirming that compression introduces a performance trade-off. Notably, $M = 2$ performs
817 the worst across all corpora, exhibiting the highest perplexity values. For merge sizes $M = 3$, $M = 4$, and $M = 5$,
818 perplexity scores are nearly identical, suggesting that moderate compression can be achieved without significantly sacrificing
819
820
821
822
823
824

language modeling quality—provided $M = 2$ is avoided. This consistency across loss and perplexity metrics further supports the robustness of maximum merge size $M = 3$ as the most effective trade-off point.

A.2. Hyper-encoder architecture

We ablate the architecture of the hyper-encoder to evaluate its effect on language modeling performance, as shown in Table 8. We compare increasingly expressive architectures, starting with a simple averaging method that introduces no additional parameters. This baseline yields the highest perplexity, highlighting its limited capacity. Adding a single attention layer significantly improves performance, and further gains are observed with a 1-layer transformer encoder. The 2-layer transformer offers marginal additional benefit, suggesting that a lightweight transformer (1–2 layers) is sufficient for effective hyper-token modeling.

Figure 6 illustrates the effect of hyper-encoder architecture on zip2zip training loss. We observe that the simple averaging method converges the fastest but plateaus at a relatively high loss, reflecting its limited capacity. As model complexity increases—with attention and transformer layers—the convergence becomes slower, yet the final loss is significantly lower. Notably, the 1-layer and 2-layer transformer encoders yield the best performance, demonstrating that additional parameters enable the model to better capture structure, albeit at the cost of slower training dynamics.

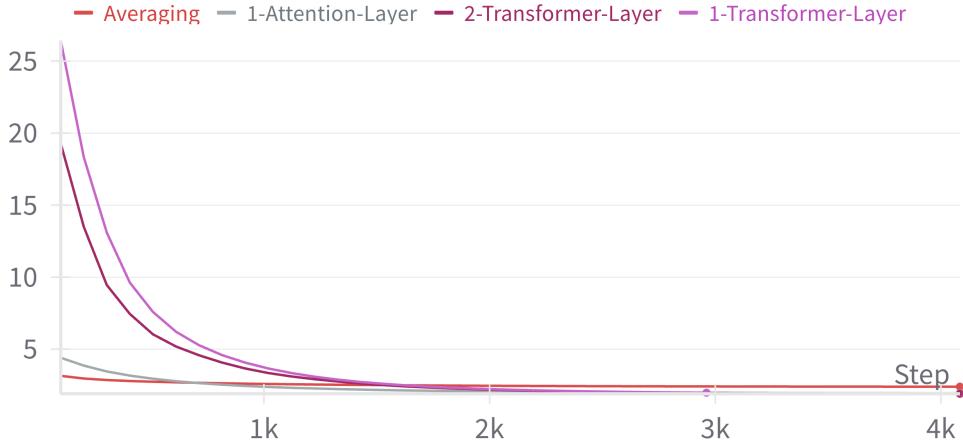


Figure 6: Effect of hyper-encoder architecture on zip2zip training loss. Averaging (no additional parameters) converges quickly but to a higher loss. As architectural complexity increases—from attention to transformer layers—convergence becomes slower, but the final loss is lower. This highlights a trade-off between training speed and modeling capacity.

B. FLOPs Estimation for zip2zip

Following the assumptions of Kaplan et al. (2020), we estimate training FLOPs (Γ) as:

$$\Gamma \approx 6 \cdot N_{\text{tokens}} \cdot N_{\text{params}},$$

where N_{tokens} is the total number of processed tokens and N_{params} is the number of trainable parameters. This estimate ignores the quadratic attention cost, assuming:

$$12 \cdot d_{\text{model}} \ll \text{sequence length}.$$

For zip2zip, this becomes:

Table 8: Ablation of hyper-encoder architecture on byte-perplexity (\downarrow) across four corpora using a 1024-token context window. Performance improves with increasingly expressive architectures.

Model	Method	Wiki	Pile	mC4	dC4
Phi-3.5-4B	averaging	1.81	1.97	2.29	2.08
	1-attention-layer	1.73	1.86	2.16	2.01
	1-transformer-layer	1.71	1.83	2.13	1.99
	2-transformer-layer	1.72	1.84	2.15	2.00

880
881 $\Gamma_{z2z} \approx 6 \cdot N_{\text{tokens}} \cdot \rho \cdot N_{\text{params}}(1 + \alpha)$,
882

883 where ρ is the compression ratio, and α accounts for the overhead of the hyper-encoder applied at the embedding and LM
884 head. The relative FLOPs ratio is then:

885
886
$$\frac{\Gamma_{z2z}}{\Gamma} = \rho \cdot (1 + \alpha).$$

887

888 Assuming the hyper-encoder mirrors the base model’s configuration, we estimate:
889

890
891
$$\alpha \approx \frac{lM}{L},$$

892

893 where l is the number of hyper-encoder layers, M is the maximum merge size, and L is the number of base model layers.
894 We illustrate this estimate across several model scales in Table 9, showing that the relative FLOPs overhead from the
895 hyper-module remains modest (typically under 15%).
896

Model	L	M	l	$\alpha = \frac{lM}{L}$
LLM-4B	14	2	1	0.14
LLM-7B	32	2	2	0.13
LLM-70B	80	3	3	0.11
LLM-400B	128	3	4	0.09

903 Table 9: Relative FLOPs overhead from the hyper-module across different model sizes.
904

905 C. Additional Results

906 Machine Translation

907 We report standard deviations for machine translation results across WMT benchmarks in Table 10, computed using the
908 lm-evaluation-harness codebase.
909

910 Table 10: Standard deviations (bootstrapped) for machine translation scores across WMT benchmarks.
911

Model	Method	WMT14 En-Fr			WMT16 En-De			WMT16 En-Ro		
		BLEU	CHRF	TER	BLEU	CHRF	TER	BLEU	CHRF	TER
Phi-3-4B	+Raw	± 2.1	± 1.4	± 1.7	± 1.9	± 1.6	± 1.8	± 1.5	± 1.3	± 2.4
	+Finetune	± 2.2	± 1.6	± 1.8	± 1.8	± 1.4	± 1.7	± 1.4	± 1.5	± 2.3
	+zip2zip	± 1.9	± 1.5	± 2.0	± 1.7	± 1.6	± 1.9	± 1.6	± 1.4	± 2.5
Phi-3-14B	+Raw	± 2.0	± 1.4	± 1.9	± 2.0	± 1.3	± 1.7	± 1.5	± 1.4	± 2.2
	+Finetune	± 2.2	± 1.4	± 1.9	± 2.0	± 1.3	± 1.9	± 1.4	± 1.3	± 2.9
	+zip2zip	± 2.1	± 1.5	± 1.8	± 2.1	± 1.6	± 1.8	± 1.5	± 1.3	± 2.6

923 D. Technical Details

924 Model and Training Configuration

- 925 • **Pretrained Model:** microsoft/Phi-3-medium-4k-instruct
926 • **Sequence Length:** 1024
927 • **Total Batch Size:** 32,768 tokens
928 • **Learning Rate Schedule:** Cosine decay
929 • **Learning Rate Range:** Max = 3e-4, Min = 1e-5
930

-
- 935 • **LoRA rank and alpha value:** Both are 32
936 • **Training Steps:** 10,000
937 • **Validation Interval:** Every 100 steps
938 • **Checkpoint Interval:** Every 500 steps
939 • **Pytorch Model Compilation:** Enabled
940
941

942 **LoRA Configuration**

- 943 • **Rank:** 16
944 • **Alpha:** 16
945 • **Target Modules:** qkv_proj, o_proj, gate_proj, down_proj, up_proj
946
947

948 **System and Libraries**

- 949 • **Hardware:** 4 × NVIDIA A100-SXM4-80GB GPUs, 64-core CPU (128 threads)
950 • **Key Libraries:**
951 – PyTorch >= 2.5.0
952 – Transformers >= 4.47.0
953 – Datasets <= 3.1.0
954 – Accelerate >= 0.26.0
955
956

957 **Compute Resources**

958 We report the compute resources used for training our models in Table 11. All training was conducted on internal servers equipped with NVIDIA H100 GPUs. We estimate GPU-hours by multiplying wall-clock training time by the number of GPUs used. No additional compute was used beyond the reported experiments; we did not perform parameter grid search, large-scale hyperparameter tuning, or exploratory runs that were excluded from the paper.

959 Table 11: Training compute resources for zip2zip experiments.
960

Model	GPUs	Time	GPU Type	GPU-Hours
Phi-3.5-Medium (14B)	4	15h 46m	NVIDIA H100 80GB	63.0
Phi-3.5-Mini (4B)	2	7h 0m	NVIDIA H100 80GB	14.0

961 **Inference.** All evaluations complete within 1 hour on a single A100 GPU, demonstrating the runtime efficiency of
962 zip2zip.
963

964 **E. Data Mixture**

965 To support effective fine-tuning, we construct a curated dataset with balanced representation across diverse domains,
966 including code, mathematics, dialogue, general web content, and multilingual text. The final dataset contains approximately
967 1 billion compressed tokens.

968 Table 12 summarizes the constituent datasets and their respective proportions. A visualization of the dataset composition
969 and sequence length characteristics is shown in Figure 7.

970 The multilingual subset in fineweb-2 includes the following languages: Mandarin Chinese (cmn_Hani), German (deu_Latn),
971 Japanese (jpn_Jpan), Spanish (spa_Latn), French (fra_Latn), Italian (ita_Latn), Portuguese (por_Latn), Dutch
972 (nld_Latn), and Arabic (arb_Arab).

973 **F. Token Stream Visualization**

Dataset	Domain	Proportion (%)
HuggingFaceFW/fineweb-edu(Lozhkov et al., 2024a)	Web / Knowledge	20%
devngho/the-stack-llm-annotations-v2(Lozhkov et al., 2024b)	Code	25%
AI-MO/NuminaMath-1.5(LI et al., 2024)	Math	20%
HuggingFaceH4/ultrachat_200k(Ding et al., 2023)	Chat / Dialogue	20%
HuggingFaceFW/fineweb-2(Penedo et al., 2024)	Multilingual	15%

Table 12: Training data composition across domains.

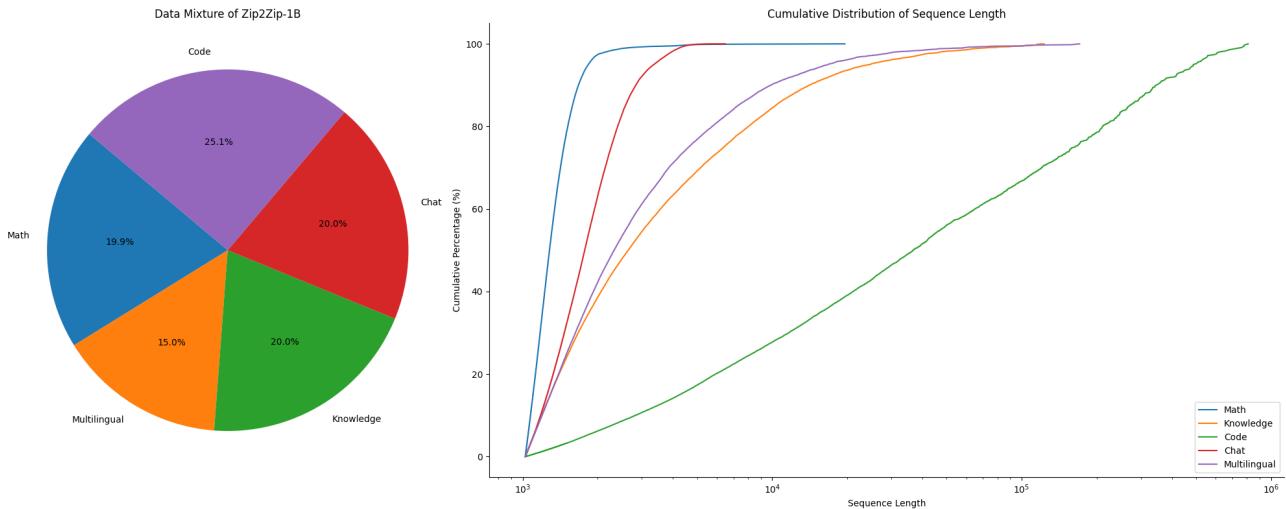


Figure 7: Left: Proportional breakdown of the fine-tuning dataset across five domains. Right: Cumulative distribution of input sequence lengths per domain (log scale). Code and multilingual data exhibit longer tail distributions, indicating greater variability in sequence lengths.

```

1045 class GPT2Attention(nn.Module):
1046     def __init__(self, config, is_cross_attention=False, layer_idx=None):
1047         super().__init__()
1048         self.config = config
1049         max_positions = config.max_position_embeddings
1050         self.register_buffer('
1051             bias',
1052             torch.tril(torch.ones(max_positions, max_positions), dtype=torch.bool).view(
1053                 1, 1, max_positions, max_positions
1054             ),
1055             persistent=False)
1056             No.2
1057             self.register_buffer('masked_bias', torch.tensor(-1e4), persistent=False)
1058             self.embed_dim = config.hidden_size
1059             self.num_heads = config.num_attention_heads
1060             self.head_dim = self.embed_dim // self.num_heads
1061             self.split_size = self.embed_dim // self.embed_dim
1062             No.3
1063             self.scale_attn_weights = config.scale_attn_weights
1064             self.is_cross_attention = is_cross_attention
1065             No.4
1066             # Layer-wise attention scaling, reordering, and upcasting
1067             self.scale_attn_by_inverse_layer_idx = config.scale_attn_by_inverse_layer_idx
1068             self.layer_idx = layer_idx
1069             self.reorder_and_upcast_attn = config.reorder_and_upcast_attn
1070             if self.is_cross_attention:
1071                 self.c_attn = Conv1D2 * self.embed_dim, self.embed_dim)
1072                 self.o_attn = Conv1D2 * self.embed_dim, self.embed_dim)
1073             else:
1074                 self.c_attn = Conv1D3 * self.embed_dim, self.embed_dim)
1075                 self.c_proj = Conv1D(self.embed_dim, self.embed_dim)
1076                 self.attn_dropout = nn.Dropout(config.attn_pdrop)
1077                 self.resid_dropout = nn.Dropout(config.resid_pdrop)
1078                 self.is_causal = True
1079                 self.pruned_heads = set()

```

(a) Default Tokenization of some Python code.

```

class GPT2Attention(nn.Module):
    def __init__(self, config, is_cross_attention=False, layer_idx=None):
        super().__init__()
        self.config = config
        max_positions = config.max_position_embeddings
        self.register_buffer('
            bias',
            torch.tril(torch.ones(max_positions, max_positions), dtype=torch.bool).view(
                1, 1, max_positions, max_positions
            ),
            persistent=False,
            No.2
            self.register_buffer('masked_bias', torch.tensor(-1e4), persistent=False)
            self.embed_dim = config.hidden_size
            self.num_heads = config.num_attention_heads
            self.head_dim = self.embed_dim // self.num_heads
            self.split_size = self.embed_dim // self.embed_dim
            No.3
            self.scale_attn_weights = config.scale_attn_weights
            self.is_cross_attention = is_cross_attention
            No.4
            # Layer-wise attention scaling, reordering, and upcasting
            self.scale_attn_by_inverse_layer_idx = config.scale_attn_by_inverse_layer_idx
            self.layer_idx = layer_idx
            self.reorder_and_upcast_attn = config.reorder_and_upcast_attn
            if self.is_cross_attention:
                self.c_attn = Conv1D2 * self.embed_dim, self.embed_dim)
                self.o_attn = Conv1D2 * self.embed_dim, self.embed_dim)
            else:
                self.c_attn = Conv1D3 * self.embed_dim, self.embed_dim)
                self.c_proj = Conv1D(self.embed_dim, self.embed_dim)
                self.attn_dropout = nn.Dropout(config.attn_pdrop)
                self.resid_dropout = nn.Dropout(config.resid_pdrop)
                self.is_causal = True
                self.pruned_heads = set()

```

(b) The same code with adaptive tokenization.

No.1
Main article: Transcription (genetics)
Transcription is when RNA is copied from DNA. During transcription, RNA polymerase makes a copy of a gene from the DNA to mRNA as needed. This process differs slightly in eukaryotes and prokaryotes. One notable difference is that prokaryotic RNA polymerase associates with DNA-processing enzymes during transcription so that processing can proceed during transcription. Therefore, this causes the new mRNA strand to become double stranded by producing a complementary strand known as the tRNA strand, which when combined are unable to form structures from base-pairing. Moreover, the template for tRNA A is the complementary strand of tRNA, which is identical in sequence to the anticodon sequence that the DNA binds to. The short-lived, unprocessed or partially processed product is termed precursor mRNA, or pre-mRNA; once completely processed, it is termed mature mRNA. [citation needed]

No.2 Uracil substitution for thymine

No.3
Main article: Post-transcriptional modification
Uracil is used instead of thymine (T) in DNA; uracil (U) is the complementary base to adenine (A) during transcription instead of thymine (T). Thus, when using a template strand of DNA to build RNA, thymine is replaced with uracil. This substitution allows the mRNA to carry the appropriate genetic information from DNA to the ribosome for translation. Regarding the natural history, uracil came first than thymine; evidence suggests that RNA came before DNA in evolution.[1] The RNA World hypothesis proposes that life began with RNA molecules, before the emergence of DNA genomes and coded proteins. In DNA, the evolutionary substitution of thymine for uracil may have increased DNA stability and improved the efficiency of DNA replication.[2][3]

No.4 Eukaryotic pre-mRNA processing

No.5
Main article: Post-transcriptional modification
DNA gene is transcribed to pre-mRNA, which is then processed to form a mature mRNA, and then lastly translated by a ribosome to a protein.

No.6 Eukaryotic pre-mRNA processing

Processing of mRNA differs greatly among eukaryotes, bacteria, and archaea. Non-eukaryotic mRNA is, in essence, mature upon transcription and requires no processing, except in rare cases.[4] Eukaryotic pre-mRNA, however, requires several processing steps before its transport to the cytoplasm and its translation by the ribosome.

No.7 Eukaryotic pre-mRNA processing

Main article: Post-transcriptional modification

No.1
Main article: Transcription (genetics)
Transcription is when RNA is copied from DNA. During transcription, RNA polymerase makes a copy of a gene from the DNA to mRNA as needed. This process differs slightly in eukaryotes and prokaryotes. One notable difference is that prokaryotic RNA polymerase associates with DNA-processing enzymes during transcription so that processing can proceed during transcription. Therefore, this causes the new mRNA strand to become double stranded by producing a complementary strand known as the tRNA strand, which when combined are unable to form structures from base-pairing. Moreover, the template for tRNA A is the complementary strand of tRNA, which is identical in sequence to the anticodon sequence that the DNA binds to. The short-lived, unprocessed or partially processed product is termed precursor mRNA, or pre-mRNA; once completely processed, it is termed mature mRNA. [citation needed]

No.2 Uracil substitution for thymine

No.3
Main article: Post-transcriptional modification
Uracil is used instead of thymine (T) in DNA; uracil (U) is the complementary base to adenine (A) during transcription instead of thymine (T). Thus, when using a template strand of DNA to build RNA, thymine is replaced with uracil. This substitution allows the mRNA to carry the appropriate genetic information from DNA to the ribosome for translation. Regarding the natural history, uracil came first than thymine; evidence suggests that RNA came before DNA in evolution.[1] The RNA World hypothesis proposes that life began with RNA molecules, before the emergence of DNA genomes and coded proteins. In DNA, the evolutionary substitution of thymine for uracil may have increased DNA stability and improved the efficiency of DNA replication.[2][3]

No.4 Eukaryotic pre-mRNA processing

No.5
Main article: Post-transcriptional modification
DNA gene is transcribed to pre-mRNA, which is then processed to form a mature mRNA, and then lastly translated by a ribosome to a protein.

No.6 Eukaryotic pre-mRNA processing

Processing of mRNA differs greatly among eukaryotes, bacteria, and archaea. Non-eukaryotic mRNA is, in essence, mature upon transcription and requires no processing, except in rare cases.[4] Eukaryotic pre-mRNA, however, requires several processing steps before its transport to the cytoplasm and its translation by the ribosome.

No.7 Eukaryotic pre-mRNA processing

Main article: Post-transcriptional modification

No.1
L'impressionnisme est un mouvement pictural français en France dans les années 1860 en opposition à l'art académique et visant à représenter le caractère éphémère de la lumière et ses effets sur les lieux et les formes. Le mouvement des impressionnistes se forme autour d'Edouard Manet, chef de file de l'avant-garde impressionniste dans les années 1860, qui ne participe à aucune exposition impressionniste. Après plusieurs scandales et refus au Salon, la grande exposition annuelle d'artistes organisée par l'Académie des Beaux-Arts, de jeunes artistes décident de s'associer pour organiser des expositions indépendantes. Cette idée se concrétise en 1874, dans une exposition qui réunit trente artistes dont Paul Cézanne, Edgar Degas, Claude Monet, Berthe Morisot, Camille Pissarro, Auguste Renoir et Alfred Sisley. Le journaliste satirique Louis Leroy invente alors le terme « impressionnisme » à partir du tableau Impression, soleil levant de Manet, devient depuis le nom du mouvement. Les artistes subissent d'abord des critiques violentes de la part de la presse et du public, mais ils sont soutenus par des collectionneurs qui permettent la tenue de leurs premières expositions, notamment Gustave Caillebotte.

No.2 L'impressionnisme commence à être accepté en 1886, grâce au soutien du nouveau gouvernement de Léon Gambetta et de critiques comme Émile Zola. Les œuvres font partie de leur entrée dans les musées, ou sont achetées par des artistes français, qui succèdent au Salon de l'Académie des Beaux-Arts et sur le marché de l'art.

No.3
Le mouvement Peinture à l'huile devient plus crucial dans le succès et la diffusion de l'impressionnisme, qui s'exporte aux États-Unis à partir de 1886, grâce à la peintre Mary Cassatt. Le mouvement y obtient un grand succès, qui participe à la consécration de Manet et au développement d'écoles impressionnistes hors de France au cours des années 1890. Cette décadence voit la mort de Morisot, Caillebotte et Sisley et la dispersion du groupe, tandis que se développent de nouvelles avant-gardes auxquelles adhèrent certains impressionnistes, comme Cézanne et Pissarro.

No.4 Les artistes impressionnistes créent une nouvelle esthétique opposée à l'art académique. Leur style apparaît pour la première fois dans les toiles peintes par Monet et Renoir à l'île de la Grenouillère, en 1869. Ils font primer la couleur sur le dessin, utilisent des compositions inhabituelles et une touche rapide, et peignent généralement en plein air sur le motif. Tournés vers des sujets modernes, ils représentent principalement des paysages, des scènes de la vie intime et les loisirs de leur époque.

No.5
L'impressionnisme dans les collections muséales
L'impressionnisme étant un mouvement largement porté par des artistes français, il est logique de trouver de nombreuses œuvres dans des musées situés en France. Toutefois, la plupart des grandes collections d'art moderne à travers le monde s'efforcent également de présenter au moins quelques exemples de toiles impressionnistes.

No.6 Les impressionnistes dans les collections muséales

No.1
L'impressionnisme est un mouvement pictural français en France dans les années 1860 en opposition à l'art académique et visant à représenter le caractère éphémère de la lumière et ses effets sur les lieux et les formes. Le mouvement des impressionnistes se forme autour d'Edouard Manet, chef de file de l'avant-garde impressionniste dans les années 1860, qui ne participe à aucune exposition impressionniste. Après plusieurs scandales et refus au Salon, la grande exposition annuelle d'artistes organisée par l'Académie des Beaux-Arts, de jeunes artistes décident de s'associer pour organiser des expositions indépendantes. Cette idée se concrétise en 1874, dans une exposition qui réunit trente artistes dont Paul Cézanne, Edgar Degas, Claude Monet, Berthe Morisot, Camille Pissarro, Auguste Renoir et Alfred Sisley, le journaliste satirique Louis Leroy invente alors le terme « impressionnisme » à partir du tableau Impression, soleil levant de Manet, devient depuis le nom du mouvement. Les artistes subissent d'abord des critiques violentes de la part de la presse et du public, mais ils sont soutenus par des collectionneurs qui permettent la tenue de leurs premières expositions, notamment Gustave Caillebotte.

No.2 L'impressionnisme commence à être accepté en 1886, grâce au soutien du nouveau gouvernement de Léon Gambetta et de critiques comme Émile Zola. Les œuvres font partie de leur entrée dans les musées, ou sont achetées par des artistes français, qui succèdent au Salon de l'Académie des Beaux-Arts et sur le marché de l'art.

No.3
Le mouvement Peinture à l'huile devient plus crucial dans le succès et la diffusion de l'impressionnisme, qui s'exporte aux États-Unis à partir de 1886, grâce à la peintre Mary Cassatt. Le mouvement y obtient un grand succès, qui participe à la consécration de Manet et au développement d'écoles impressionnistes hors de France au cours des années 1890. Cette décadence voit la mort de Morisot, Caillebotte et Sisley et la dispersion du groupe, tandis que se développent de nouvelles avant-gardes auxquelles adhèrent certains impressionnistes, comme Cézanne et Pissarro.

No.4 Les artistes impressionnistes créent une nouvelle esthétique opposée à l'art académique. Leur style apparaît pour la première fois dans les toiles peintes par Monet et Renoir à l'île de la Grenouillère, en 1869. Ils font primer la couleur sur le dessin, utilisent des compositions inhabituelles et une touche rapide, et peignent généralement en plein air sur le motif.

No.5
Les impressionnistes dans les collections muséales
L'impressionnisme étant un mouvement largement porté par des artistes français, il est logique de trouver de nombreuses œuvres dans des musées situés en France. Toutefois, la plupart des grandes collections d'art moderne à travers le monde s'efforcent également de présenter au moins quelques exemplaires de toiles impressionnistes.

No.6 Les impressionnistes dans les collections muséales

(c) Default Tokenization of some biomedical text.

(d) The same text with adaptive tokenization.

(e) Default Tokenization of text in French.

(f) The same text with adaptive tokenization.