# BITUNE: LEVERAGING BIDIRECTIONAL ATTENTION TO IMPROVE DECODER-ONLY LLMS

Anonymous authors

004

010 011

012

013

014

015

016

017

018

019

021

Paper under double-blind review

#### ABSTRACT

Decoder-only large language models typically rely solely on masked causal attention, which limits their expressiveness by restricting information flow to one direction. We propose Bitune, a method that enhances pretrained decoder-only LLMs by incorporating bidirectional attention into prompt processing. We evaluate Bitune in instruction-tuning and question-answering settings, showing significant improvements in performance on commonsense reasoning, arithmetic, and language understanding tasks. Furthermore, extensive ablation studies validate the role of each component of the method, and demonstrate that Bitune is compatible with various parameter-efficient finetuning techniques and full model finetuning.

022 1 INTRODUCTION

Large Language Models (LLMs) are being deployed in numerous practical applications where
humans engage with them through various forms of natural language interaction. In use cases such as
general purpose assistants (OpenAI, 2024), medical diagnosticians (Thirunavukarasu et al., 2023),
game-conversation generation (Cox & Ooi, 2023) or coding-assistants (Roziere et al., 2023), the
ability for an LLM to precisely interpret and respond to user inputs is of primary concern.

Correspondingly, Instruction-Tuning (IT) (Chung et al., 2024; Ouyang et al., 2022a) is the prevailing paradigm for finetuning LLMs after their self-supervised pretraining phase to improve them for such tasks. Here, the model is trained on a dataset comprised of pairs of instructions and corresponding responses. Given the instruction-with-response structure of IT data, the generation of an LLM response can be divided into two phases: first, converting the instruction into *key* and *value* embeddings, which we refer to as instruction features; second, using these features to autoregressively generate an answer. Due to this task's inherently conditional nature, the instruction features' effectiveness is crucial for obtaining high-quality model outputs.

In the past, bidirectional attention (Schuster & Paliwal, 1997) has been a key technique for obtaining
stronger features for words or tokens. This is because the meaning of a word depends greatly on
its context. In particular, for some words in a sentence, the information that comes later might be
far more informative for generating a meaningful representation and resolving ambiguities. With
only uni-directional causal attention, where the representation of each word is restricted to depend
solely on the words that came before, this cannot be achieved. This is the reason why many previous
transformers such as encoder-only BERT (Devlin et al., 2019) and encoder-decoder T5 (Raffel et al.,
2020) employed bidirectional attention to improve the encoding of the input and why tasks like text
retrieval (Lewis et al., 2020; Li & Li, 2023) still rely on this.

However, in the context of LLMs, architectures utilizing bidirectional attention have fallen out of favor, as decoder-only models such as GPT (OpenAI, 2024) and Llama (AI@Meta, 2024) have focused on and vastly improved the generative performance of language models. These architectures are trained by large volumes of data with next-token prediction, eschewing any look-ahead mechanism for the sake of better autoregressive modeling. As there is simply more unlabeled data available for pretraining, training a decoder-only architecture on unlabeled data, and then finetuning it for tasks with instruction-tuning, is the best modus operandi of today (Wang et al., 2022). However, with this switch to decoder-only architectures, we lost bidirectional attention in the process. As we know this can improve feature representations for instructions, we set out to re-introduce bidirectional attention, such that it can be integrated into pretrained decoder-only LLMs.



Figure 1: **Overview of Bitune.** (a) During the prefilling phase, features are obtained from the prompt using both causal and bidirectional attention in two passes with separate weights. The two sets of keys and values are then combined using a weighted average before being passed to the decoding phase. (b) During the decoding phase, new tokens are generated in the standard way with causal attention, utilizing the features extracted from the instruction in the previous step, along with the features of other generated tokens.

Our new method Bitune adds bidirectional attention to decoder-only architectures and combines it with causal attention to generate two sets of instruction features, using two different sets of weights. These features are then integrated, utilizing learnable mixing coefficients, and later used as the KV-cache for response generation. Notably, the autoregressive response generation process remains unaffected by the bidirectional attention and continues to be causal. By realizing these adaptations with parameter-efficient finetuning methods, we introduce only a minimal set of new parameters.

Overall, our contributions are as follows:

066

067

068

069

071

079

081

082

084

090

096

- We propose a novel method, Bitune, that improves the performance of pretrained decoderonly LLMs in instruction-following and question-answering settings.
- We evaluate the method on multiple downstream tasks, showing consistent improvements over the baselines.
- We conduct an extensive ablation study investigating the necessity of each component of the method, and showing the method's PEFT-agnosticism, as well as its effectiveness in full finetuning scenarios.

#### 2 **BIDIRECTIONAL INSTRUCTION-TUNING**

r

In the instruction-tuning setting (Ouyang et al., 2022b; Zhang et al., 2024), a dataset  $\mathcal{D}$  consists of instruction-answer pairs that are used to adapt the model in a supervised fashion. Formally, a dataset of size N can be described as  $\mathcal{D} = \{q, a\}_{i=1}^{N}$ , where q and a are instructions and answers. The training objective is to model p(a|q) in an autoregressive manner: This means the answer is generated one token at a time, such that token  $a_i$  at position *i* has access to all earlier tokens:

$$p(a|q) = \prod_{i=1}^{|a|} p(a_i|a_1, \dots a_{i-1}, q),$$

where |a| denotes the length of the answer. Note how compared to the regular language modeling objective, the response is already conditional (on the instruction q) even for the first generated token.

This naturally leads response-generation to be divided into two phases: prefilling and decoding. During the prefilling phase, the entire instruction – also often called a prompt – is processed concurrently to generate a series of features to be stored. For a Transformer architecture (Vaswani et al., 2017), these features are those of the key and value vectors, which can be stored in a KV-cache to avoid costly recomputations. During the subsequent decoding phase, the model generates output tokens sequentially, one token at a time, based on the KV-cache of the instruction and the already generated tokens.

107 In this work, we introduce Bitune, a method to leverage this two-phase process to improve instructiontuning of language models. In our approach, the model processes the instruction with both causal

108

#### Algorithm 1 Python-like pseudocode of Bitune inference.

```
110
      # prompt - tensor with tokenized instruction
      # theta - tensor with mixing coeff. for each layer
111
      # theta_init - initial value of mixing coefficients
112
      k_c, v_c = model_causal(prompt) # Extract causal features
113
      k_b, v_b = model_bidir(prompt) # Extract bidirectional features
114
      # Combine both sets of features
115
      alpha = theta.abs() / (theta.abs() + theta_init)
116
      k = k_c * (1 - alpha) + k_b * alpha
      v = v_c * (1 - alpha) + v_b * alpha
117
      kv = (k, v)
118
119
      c_token = SEP # Initialize generation with a predefined token
120
      answer = [c_token]
      while c_token != EOS: # Stop generation at the end-of-sequence token
121
         # Get features of current token and logits of next token
122
         k, v, logits = model_causal(c_token, kv)
         kv = concat(kv, (k, v)) # Concatenate it with current KV cache
123
         c_token = get_token(logits) # Determine next predicted token
124
         answer.append(c_token) # Append generated token to the answer
125
```

126 127

and bidirectional attention using separate sets of parameters, leading to an enhanced KV-cache that is
 then used to condition the answer. Figure 1 provides an overview of the method, while Algorithm 1
 presents pseudocode for the inference process.

Two Sets of Features. In Bitune, the model performs two passes on the instruction to obtain two
 kinds of features for every transformer block. Namely, a set of causal features that the model was
 originally trained to process and utilize,

$$K_{c} = X_{c} W_{kc}, \quad V_{c} = X_{c} W_{vc}, \tag{1}$$

and a set of bidirectional features encoding the instruction without the constraints of causal masking,

134 135 136

$$K_b = X_b W_{kb}, \quad V_b = X_b W_{vb}. \tag{2}$$

To allow the model to learn how to process the causal and bidirectional features differently, we introduce two sets of weights: one for the bidirectional pass on the instruction  $(W_{kb}, W_{vb})$  and another for the causal pass on the instruction, which is also used for the causal generation of answer tokens  $(W_{kc}, W_{vc})$ .

In the case of the first block of the model, representations  $X_c$ ,  $X_b$  are the initial token embeddings. In other cases, they are the output of the preceding block and were processed by different components including the self-attention mechanism, which can be defined as:

148 149

152

160 161

Attention
$$(Q, K, V, M) = \operatorname{softmax}\left(\frac{QK^T}{\sqrt{d_k}} + M\right)V,$$
 (3)

where Q are the queries, M is the attention mask, and  $d_k$  is the dimension of keys and queries.

For the causal pass, the mask  $M_c$  enforces causality by masking future tokens, such that tokens j can only attend to earlier tokens  $i \le j$ , while for the bidirectional pass, no masking is applied:

$$M_c(i,j) = \begin{cases} 0 & \text{if } i \le j \\ -\infty & \text{if } i > j \end{cases}, \qquad M_b(i,j) = 0.$$

$$\tag{4}$$

The final KV-cache is obtained by a learnable convex combination of causal and bidirectional features,

$$K_{\text{Bitune}} = K_{\boldsymbol{c}} \cdot (1 - \alpha_k) + K_{\boldsymbol{b}} \cdot \alpha_k, \quad V_{\text{Bitune}} = V_{\boldsymbol{c}} \cdot (1 - \alpha_v) + V_{\boldsymbol{b}} \cdot \alpha_v, \tag{5}$$

Table 1: Zero-shot results after instruction-tuning on the UltraFeedback dataset. We compare Bitune to the performance of the original model, the model finetuned with LoRA, and with LoRA<sub>16</sub> using two times higher rank to match the parameter count of our method. Bitune significantly outperforms the baselines on almost all tasks for all models. 

		PIQA	ARC	CSQA	SIQA	MMLU	Avg.
Model	Method	-		-	-		-
Gemma-2B	Pretrained	57.5	36.9	35.5	38.2	34.0	40.4
	LoRA	66.7	43.4	42.3	44.3	31.7	45.7
	LoRA <sub>16</sub>	66.5	42.7	42.3	43.8	31.6	45.4
	Bitune	69.6	47.5	46.9	49.5	35.3	49.7
Gemma-7B	Pretrained	73.1	78.3	62.0	64.7	59.0	67.4
	LoRA	84.2	79.2	68.5	71.9	55.3	71.8
	LoRA <sub>16</sub>	83.9	79.2	68.4	72.0	53.4	71.4
	Bitune	83.6	80.1	69.2	72.7	53.8	71.9
Llama2-7B	Pretrained	59.2	38.1	32.6	45.1	36.0	42.2
	LoRA	69.5	49.9	45.3	57.0	41.1	52.6
	LoRA <sub>16</sub>	69.9	49.9	45.6	56.7	41.2	52.6
	Bitune	70.0	51.1	48.1	59.1	41.9	54.0
Llama3-8B	Pretrained	69.0	73.6	65.4	56.8	56.0	64.2
	LoRA	81.9	74.5	69.2	69.0	57.6	70.4
	LoRA <sub>16</sub>	82.4	74.9	70.5	68.6	58.0	70.9
	Bitune	84.4	77.4	72.7	70.1	59.0	72.7
Phi-2	Pretrained	70.3	67.3	61.4	65.0	45.4	61.9
	LoRA	76.3	66.7	61.6	66.6	48.2	63.9
	LoRA <sub>16</sub>	76.1	66.6	61.6	66.8	47.7	63.8
	Bitune	76.5	67.2	63.0	68.5	48.9	64.8

where  $\alpha$  represents the *bidirectional-to-causal* ratio of features. This ratio is parameterised as

 $\alpha_j = \frac{|\theta_j|}{\theta_{\mathrm{init}} + |\theta_j|}, \quad j \in \{k, v\}$ where  $\theta_i$  is a learnable *mixing coefficient* per transformer block, and  $\theta_{init}$  is a hyperparameter defining the initial value of  $\theta_i$ . The mixing coefficients are learnable to allow each block to independently adjust the balance between bidirectional and causal features throughout the training.

(6)

**Parameter Efficient Fine-tuning** Note that the components of the model, other than the key and value projections, can have their own separate sets of weights as well. In the case of full finetuning, this approach would require an additional set of full weights, which is impractical for large models.

Instead, we adapt our model using parameter-efficient finetuning methods. These introduce only a fraction of trainable parameters, making it viable to have two modified variants of the model within a single forward pass. In the default configuration of our method, we utilize the Low-Rank Adaptation (LoRA) of Hu et al. (2022) to adapt the model. However, Bitune can utilize different methods for updating the weights, including full model finetuning and other parameter-efficient techniques, as demonstrated in our ablations section.

#### EXPERIMENTS

3.1 INSTRUCTION-TUNING 

Our core experiments involve training pretrained language models on an instruction-tuning dataset and zero-shot evaluating them on downstream tasks. We evaluate Bitune on multiple models, comparing results to standard finetuning with LoRA, and zero-shot results of pretrained models without finetuning. 

Specifically, we use a subset of the cleaned UltraFeedback (Cui et al., 2023) dataset, which contains instructions and corresponding answers generated by various LLMs. From this dataset, we select completions generated by GPT-4 (OpenAI, 2024), ensuring high-quality responses for training. To fit every model on a single GPU, we filter out samples longer than 512 tokens, which leaves us with

roughly 10,000 samples for training. For results on another instruction-tuning dataset, please see the
 Appendix 6.5.

We test the method on pretrained decoder-only language models of two different scales of approximately 2 billion and 7 billion parameters. The specific models used in our experiments are: Gemma 2B and 7B (Gemma Team et al., 2024), Llama2 7B (Touvron et al., 2023) and Llama3 8B (AI@Meta, 2024), and Phi-2 (Li et al., 2023), which has 2.7 billion parameters. We use HuggingFace Transformers (Wolf et al., 2020) implementation of these models.

For updating the weights we use the HuggingFace PEFT (Mangrulkar et al., 2022) implementation of LoRA, with the default rank of 8, and apply it to all linear layers of MLP and self-attention components of the model. We compare Bitune with the following three baselines: **Pretrained** initial model without any finetuning; **LoRA** - model finetuned with LoRA without Bitune-specific modifications, using rank of 8 as used in our method; and **LoRA**<sub>16</sub> - model finetuned with LoRA, using a rank of 16 to provide a fair comparison in terms of the number of parameters, as our method introduces two sets of weights.

For each model, we tune the learning rate on the LoRA baseline using steps on the approximate logarithmic scale (1e-4, 3e-4, 1e-3, 3e-3), and then apply the same rate to the other approaches. Note that this potentially puts our method at a disadvantage compared to the LoRA baseline. All hyperparameters are reported in the Appendix 6.2.

Models are evaluated zero-shot on multiple-choice tasks to assess their performance. For commonsense reasoning, we use the PIQA (Bisk et al., 2020), CommonsenseQA (Talmor et al., 2019), ARC-Challenge (Clark et al., 2018), and SIQA (Sap et al., 2019) datasets, while for language understanding, we use the MMLU (Hendrycks et al., 2021) benchmark. Each task consists of a series of questions, each with multiple choices, where only one answer is correct. As the tasks follow the question-answer pattern, they are compatible with the instruction-tuning setting.

For evaluation, we use the *Language Model Evaluation Harness* framework (Gao et al., 2023). This framework formats each question using a predefined template, tokenizes the question-choice pairs, runs them through the model, and compares the log-likelihoods of the choices to determine the selected answer. For each model and approach configuration, we conduct experiments using three different random seeds, and average the results.

Models are loaded and trained using bfloat16 precision, except for the *mixing* part, which operates in the full 32-bit floating-point format. This high level of precision for the mixing of features is important, as minor numerical inaccuracies in the learnable coefficients and intermediate results of the mixing operation may lead to significant deviations in the model's behavior.

In the decoding phase of the inference with Bitune, to initiate generation, the model requires at least a single token to obtain a set of attention *queries*, in addition to the *keys* and *values* extracted from the instruction. To facilitate this, one can introduce a new learnable *<sep>* token that would be placed at the beginning of modeled answer, or utilize an existing token. For our experiments, we opted to move the last token of the instruction template to the beginning of the modeled answer. For details on the instruction template used, please refer to the section 6.7 of the Appendix.

Results. Table 1 shows consistent and significant gains after instruction-tuning with Bitune, with the highest gains seen on the Gemma-2B model, showing a 4 percentage point (pp) improvement over the baseline LoRA and a 9.3 pp improvement over the pretrained model. For the other models, the average gains over baseline finetuning are equal to 1.8, 1.4, and 0.9 pp, for Llama3-8B, Llama2-7B, and Phi-2 respectively.

It is worth noting that the Gemma-7B model shows the lowest average improvement across all tasks, with merely 0.1 pp gain over the baseline finetuning. It is also a single case where the baseline pretrained model achieved the highest score on a task, MMLU, with degraded performance in all fine-tuning approaches. However, this is not an issue with the model's scale, as significant gains are observed with the Llama2-7B and Llama3-8B models.

265

266 3.2 DOWNSTREAM TASK TRAINING 267

This complementary experiment verifies whether Bitune increases the capacity of the model within
 the narrow scope of a single task. It follows the setup from the instruction-tuning experiments with a few changes. Namely, models are not instruction-tuned but trained separately for each evaluation task

Table 2: Result for the downstream task training. We show accuracy on downstream tasks for the
baseline LoRA finetuning and Bitune, averaged over 3 seeds. On average, our method works better
than standard LoRA. We see the most significant gains on the GSM8K dataset, but slightly lower
results for Gemma-7B and the SIQA task.

		PIQA	ARC	CSQA	SIQA	GSM8K	Avg.
Model	Method						<sup>c</sup>
Gemma-2B	LoRA	81.4	58.0	77.2	77.4	30.2	64.8
	LoRA <sub>16</sub>	81.1	59.1	77.4	77.1	30.2	65.0
	Bitune	83.3	60.0	78.3	76.6	33.0	66.2
Gemma-7B	LoRA	91.4	84.6	84.4	79.4	59.1	79.8
	LoRA <sub>16</sub>	91.6	83.9	83.9	79.7	59.4	79.7
	Bitune	92.1	84.2	84.2	79.4	59.4	79.9
Llama2-7B	LoRA	84.4	66.6	81.5	82.7	32.0	69.4
	LoRA <sub>16</sub>	84.4	66.8	81.7	82.3	31.1	69.3
	Bitune	84.4	66.9	82.0	81.4	32.9	69.5
Llama3-8B	LoRA	90.2	80.7	83.9	83.1	60.4	79.7
	LoRA <sub>16</sub>	90.4	81.3	83.4	83.1	59.6	79.6
	Bitune	90.5	81.3	84.1	82.1	63.4	80.3
Phi-2	LoRA	82.8	76.3	78.7	80.3	58.6	75.3
	LoRA <sub>16</sub>	83.1	76.1	78.6	80.6	57.5	75.2
	Bitune	83.9	77.0	79.0	80.4	59.2	75.9

289 290 291

284

287

292 293

294

using the corresponding training set. We use PIQA, ARC, CSQA, and SIQA introduced earlier, and an additional arithmetic task, GSM8K (Cobbe et al., 2021).

GSM8K differs from the other tasks, where we compare log-likelihoods of predefined answers, as it
 requires the model to generate a full answer token-by-token, including the intermediate step-by-step
 reasoning. The final answer follows a specific pattern, making it feasible to extract the answer using
 methods such as regular expressions as the model learns to adhere to this pattern during training.

Results. Table 2 presents the results, demonstrating improvements when finetuning on the downstream tasks with Bitune, similar to those seen with instruction-tuning. While there are a few cases
where the baseline finetuning achieves better results on specific tasks, when considering the average
gains, applying our method is beneficial across all models. Most importantly, on the GSM8K dataset,
we see consistent high gains, suggesting that our method improves the model's reasoning ability
in generative tasks. We present additional results on GSM8K with a 22B parameter model in the
Appendix 6.6.

Similar to the instruction-tuning results, the highest gains are observed on the Gemma-2B model, while the lowest on the Gemma-7B. This indicates that the effectiveness of our method depends on the specific model used.

309 310

311

#### 3.3 Ablations

We conduct an ablation study on Bitune using the same experimental setup as in the instructiontuning experiment. For this purpose, two models are used: Gemma-2B and Llama3-8B, representing different size scales and model families.

Component Removal To verify the necessity of each component of the method, we remove selected
 parts to answer the following questions:

317 318

319

320

• Can we simply modify the attention mask to apply bidirectional attention on the prompt, without using separate weights and mixing? - We test this simplest variant, which we refer to as **Naive Bidir**.

Do we need two sets of features? Is it sufficient to obtain bidirectional features from the prompt using different weights than those used for causal answer generation? - We remove the part responsible for generating the set of causal features, and therefore also the mixing component; we refer to this as No Mixing.

Table 3: Ablation study on components of Bitune. We report zero-shot accuracy averaged over
 PIQA, ARC, CSQA, SIQA and MMLU tasks. The components are explained in section 3.3. We see
 that all ablated variants outperform the LoRA baseline, and combining all components performs the
 best. Of note, especially bidirectional attention improves results the most.

						Avg. Acc.
Model	Method	Causal	Bidir.	Mixing	Sep. Weights	-
Gemma-2B	LoRA	$\checkmark$	-	-	-	45.7
	Naive Bidir.	-	$\checkmark$	-	-	47.9
	No Mixing	-	$\checkmark$	-	$\checkmark$	48.9
	Only Causal	$\checkmark$	-	$\checkmark$	$\checkmark$	46.9
	Shared Weights	$\checkmark$	$\checkmark$	$\checkmark$	-	47.4
	Bitune	$\checkmark$	$\checkmark$	$\checkmark$	$\checkmark$	49.7
Llama3-8B	LoRA	$\checkmark$	-	-	-	70.4
	Naive Bidir.	-	$\checkmark$	-	-	71.9
	No Mixing	-	$\checkmark$	-	$\checkmark$	71.5
	Only Causal	$\checkmark$	-	$\checkmark$	$\checkmark$	71.1
	Shared Weights	$\checkmark$	$\checkmark$	$\checkmark$	-	72.3
	Bitune	$\checkmark$	$\checkmark$	$\checkmark$	$\checkmark$	72.7

 Table 4: Combining Bitune with different PEFT methods. Performances are averaged over 3 seeds.

 We can see that our method improves results regardless of the specific PEFT method used.

		PIQA	ARC	CSQA	SIQA	MMLU	Avg.
Model	Method						
Gemma-2B	LoRA	66.7	43.4	42.3	44.3	31.7	45.7
	Bitune	69.6	47.5	46.9	49.5	35.3	49.7
	DoRA	66.7	43.6	41.9	44.7	31.9	45.8
	Bitune <sub>DoRA</sub>	69.6	47.5	46.9	49.7	35.1	49.8
	IA3	67.2	46.5	45.5	37.6	32.5	45.9
	Bitune <sub>IA3</sub>	67.5	47.3	48.9	44.3	33.6	48.3
Llama3-8B	LoRA	81.9	74.5	69.2	69.0	57.6	70.4
	Bitune	84.4	77.4	72.7	70.1	59.0	72.7
	DoRA	82.1	75.4	70.2	69.2	57.7	70.9
	Bitune <sub>DoRA</sub>	84.1	77.1	72.0	70.6	58.7	72.5
	IA3	80.9	75.5	68.3	66.4	58.7	70.0
	Bitune <sub>IA3</sub>	83.4	75.7	69.2	67.8	58.8	71.0

• Are the gains solely from mixing two sets of features generated with different weights, or is bidirectional attention necessary? - Here we keep the attention mask causal to generate both sets of features, which we refer to as **Only Causal**.

• Do we need separate weights, or can the same weights be used to generate both causal and bidirectional features? - To answer this question, we do not introduce the second set of weights and use the same LoRA for both passes on the prompt, calling it **Shared Weights**.

The results, averaged over three seeds and presented in Table 3, indicate that all variants of the method
 lead to gains over the baseline LoRA finetuning. However, the highest gains are observed in the full
 variant of Bitune, demonstrating that each component contributes to the method's effectiveness.

**Different PEFT Methods** To verify the impact of different PEFT methods on the performance of our method, we compare Bitune in combination with the following techniques: LoRA (Hu et al., 2022), that reparametrizes weight updates as a multiplication of two low-rank matrices; DoRA (Liu et al., 2024), which decomposes these weight updates into direction and magnitude; and IA3 (Liu et al., 2022), that instead rescales activations with learnable vectors.

The results are shown in Table 4. We find consistent gains across all three PEFT methods we analyze, with gains ranging from +1.6% to +4.0% for averaged accuracy. This demonstrates that Bitune is PEFT-agnostic and can be combined with existing and future innovations in PEFT methods.

Table 5: Results for full finetuning on the instruction-tuning setup. We compare our method to
 full finetuning baseline. For Bitune, we optimize two sets of full model's weights.

		PIQA	ARC	CSQA	SIQA	MMLU	Avg.
Model	Method						
Gemma-2B	Full FT	69.0	46.7	43.5	43.7	34.8	47.5
	Bitune (Full FT)	70.3	48.0	47.4	43.9	36.7	49.3

**Full Finetuning** Additionally, we test whether Bitune yields gains with full finetuning (*Full-FT*), by optimizing two sets of full model's parameters. We conduct experiments on Gemma-2B model, and compare results with standard Full-FT baseline. The results in Table 5 demonstrate that Bitune improves the model's performance even in full finetuning scenarios.

**Initialization of Mixing Coefficient** The initial value of the mixing coefficient  $\theta$  is a hyperparameter in our method. To evaluate its impact on the performance and the training dynamics of the bidirectional-to-causal ratio of features, we conduct experiments on the instruction-tuning setup with the following values: 0.1, 0.01, and 0.001.

397

398

399 400

401

402

403 404

387

388

389 390

391

392





Figure 2: Bidirectional-to-causal ratio during
training. The ratio is averaged over all layers and
shown for different initial values of mixing coefficients for Llama3-8B. The initial value impacts
the change of the ratio, with higher values slowing
it down, and lower values increasing it.

Figure 3: **Ratio across layers**. Here we show the final ratio of the model in Fig. 2 across all layers for the K and V values. The utilization of bidirectional attention is spread across all layers.

Table 6a demonstrates that the initial value of the mixing coefficient impacts the performance, with 0.01 being the most optimal value for both models, regardless of their scale. Figure 2 shows that the initial value substantially affects the rate of change of the mixing ratio, with the higher value leading to nearly no change in the ratio, while the lower value results in sharp changes at the very beginning of the training. In Figure 3, we also observe that after training, all layers utilize the bidirectional attention.

Attention Mask of Second Pass We test another option for the attention mask of the second pass on the instruction. We transpose the causal attention mask, blocking information flow from the *past* tokens, and allowing from the *future* tokens - we call it *anti-causal* attention mask.

Results shown in Table 6b indicate that the instruction has to be processed with full bidirectional attention to achieve the highest gains. Combining *causal* and *anti-causal* features independently does not lead to the same high performance.

423 424

425

4 RELATED WORK

Our approach shares similarities with the concept of "prefix language modeling", which enables a
decoder-only model to handle bidirectional context within a prefix (instruction) while maintaining
causal generation for the output sequence. The prefix-LM architecture was introduced by Liu et al.
(2018) and further explored and popularized by Raffel et al. (2020). In their work on T5, Raffel et al.
(2020) pretrained the prefix-LM architecture alongside other architectures, such as encoder-decoder
and decoder-only models, demonstrating that prefix-LM outperforms decoder-only models on both
training objectives: denoising and language modeling.

Table 6: Ablation of Bitune's attention. We vary the additional attention mask for processing the instruction besides the causal pass, and evaluate different initial values of mixing coefficient,  $\theta_{init}$ .

(a) Different init. values for mixing coefficients. (b) 2nd pass attention masks for instruction features.

Model	Init. Value	Avg. Acc.	Model	Attention Mask
Gemma-2B	0.1	49.4	Gemma-2B	Causal
	0.01	49.7		Anti-causal
	0.001	47.2		Bidirectional
Llama3-8B	0.1	72.7	Llama3-8B	Causal
	0.01	72.7		Anti-causal
	0.001	72.3		Bidirectional

443 444

434

The prefix-LM approach has been used in UniLM (Dong et al., 2019), which trains a single transformer on three types of language modeling tasks: unidirectional, bidirectional, and sequence-to-sequence prediction. UniLM employs a shared Transformer network and utilizes specific self-attention masks to control the context that predictions are conditioned on, where the sequence-to-sequence task is equivalent to the prefix-LM approach.

Additionally, UL2 (Tay et al., 2023) introduces a pretraining objective called "Mixture of Denoisers", which combines various denoising strategies, including the prefix-LM approach. Lastly, XLNet (Yang et al., 2019) also allows for non-causal word ordering by allowing random permutations to be used with a next-token prediction objective.

All these works focused on the model *pretraining*. As for the utilization of pretrained causal language 454 models, Springer et al. (2024) show in their work that simply repeating the input to these models 455 improves the quality of token embeddings for text-retrieval. This work addresses the limitation 456 that token embeddings in autoregressive models cannot contain information from tokens appearing 457 later in the input. By repeating the input twice, the early tokens are allowed to encode information 458 about later tokens, thereby improving the quality of the embeddings. Another approach, LLM2Vec 459 (BehnamGhader et al., 2024), demonstrates that pretrained causal LLMs can be effectively converted 460 to BERT-like encoders. It can be done by enabling bidirectional attention, training the model on the 461 objective of masked token prediction, and applying unsupervised contrastive learning. 462

### 5 DISCUSSION

Limitations. During standard instruction-tuning training, the instruction and the answer are processed in a single forward pass. In our method, this processing is explicitly split into phases, extracting instruction features with two passes & answer modeling, increasing both training time and memory usage. As for the inference, these two passes on the instruction, bidirectional and causal, can be processed either in parallel, impacting memory usage, or sequentially, increasing latency.

However, this is a minor limitation in the context of instruction-tuning, since typically smaller datasets are used compared to pretraining, leading to relatively short training times. Furthermore, at inference time the added latency for processing the instruction is negligible, as the bulk of compute is used for the autoregressive answer generation: *e.g.* only a 0.2s increase from 11.5s to 11.7s for a 200 token response to a 50 token instruction for Llama3-8B (see Appendix 6.4 for more results).

475
 476
 476
 476
 476
 476
 476
 477
 478
 478
 478
 478
 479
 479
 479
 479
 470
 470
 470
 471
 472
 473
 474
 474
 475
 475
 476
 477
 478
 478
 478
 479
 479
 470
 470
 470
 471
 471
 472
 473
 474
 474
 475
 475
 475
 476
 477
 478
 478
 478
 478
 480
 478
 478
 478
 478
 478
 478
 478
 478
 478
 478
 478
 478
 478
 478
 478
 480
 480
 480
 480
 480
 480
 480
 480
 480
 480
 480
 480
 480
 480
 480
 480
 480
 480
 480
 480
 480
 480
 480
 480
 480
 480
 480
 480
 480
 480
 480
 480
 480
 480
 480
 480
 480
 480
 480
 480
 480
 480
 480
 480
 480
 480
 480
 480
 480
 480
 480
 480
 480

481

463

- 482
- 483
- 484
- 485

Reproducibility Statement. We have made efforts to ensure the reproducibility of our work. The code necessary to reproduce all experiments presented in this paper is provided in the supplementary materials. Information about hyperparameters and other experimental specifications can be found in the appendix. Each experiment was conducted using a single NVIDIA A100 GPU.

#### References

490 491

492

493 494

495

496

497

498

- AI@Meta. Llama 3 model card, 2024. URL https://github.com/meta-llama/llama3/ blob/main/MODEL\_CARD.md.
- Parishad BehnamGhader, Vaibhav Adlakha, Marius Mosbach, Dzmitry Bahdanau, Nicolas Chapados, and Siva Reddy. LLM2Vec: Large language models are secretly powerful text encoders. *arXiv* preprint: arXiv:2404.05961, 2024.
- Yonatan Bisk, Rowan Zellers, Ronan Le Bras, Jianfeng Gao, and Yejin Choi. Piqa: Reasoning about
   physical commonsense in natural language. In *AAAI*, 2020.
- Hyung Won Chung, Le Hou, Shayne Longpre, Barret Zoph, Yi Tay, William Fedus, Eric Li, Xuezhi
  Wang, Mostafa Dehghani, Siddhartha Brahma, Albert Webson, Shixiang Shane Gu, Zhuyun Dai,
  Mirac Suzgun, Xinyun Chen, Aakanksha Chowdhery, Sharan Narang, Gaurav Mishra, Adams Yu,
  Vincent Zhao, Yanping Huang, Andrew Dai, Hongkun Yu, Slav Petrov, Ed H. Chi, Jeff Dean, Jacob
  Devlin, Adam Roberts, Denny Zhou, Quoc V. Le, and Jason Wei. Scaling instruction-finetuned
  language models. *JMLR*, 2024.
- Peter Clark, Isaac Cowhey, Oren Etzioni, Tushar Khot, Ashish Sabharwal, Carissa Schoenick, and Oyvind Tafjord. Think you have solved question answering? try arc, the ai2 reasoning challenge, 2018.
- Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, Mark Chen, Heewoo Jun, Lukasz Kaiser,
  Matthias Plappert, Jerry Tworek, Jacob Hilton, Reiichiro Nakano, Christopher Hesse, and John
  Schulman. Training verifiers to solve math word problems, 2021.
- Samuel Rhys Cox and Wei Tsang Ooi. Conversational interactions with npcs in llm-driven gaming:
   Guidelines from a content analysis of player feedback. In *International Workshop on Chatbot Research and Design*, 2023.
- Ganqu Cui, Lifan Yuan, Ning Ding, Guanming Yao, Wei Zhu, Yuan Ni, Guotong Xie, Zhiyuan Liu, and Maosong Sun. Ultrafeedback: Boosting language models with high-quality feedback. *arXiv:2310.01377*, 2023.
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. BERT: Pre-training of deep
   bidirectional transformers for language understanding. In ACL, 2019.
- Li Dong, Nan Yang, Wenhui Wang, Furu Wei, Xiaodong Liu, Yu Wang, Jianfeng Gao, Ming Zhou, and Hsiao-Wuen Hon. Unified language model pre-training for natural language understanding and generation. In *NeurIPS*, 2019.
- Leo Gao, Jonathan Tow, Baber Abbasi, Stella Biderman, Sid Black, Anthony DiPofi, Charles Foster, Laurence Golding, Jeffrey Hsu, Alain Le Noac'h, Haonan Li, Kyle McDonell, Niklas Muennighoff, Chris Ociepa, Jason Phang, Laria Reynolds, Hailey Schoelkopf, Aviya Skowron, Lintang Sutawika, Eric Tang, Anish Thite, Ben Wang, Kevin Wang, and Andy Zou. A framework for few-shot language model evaluation, 2023. URL https://zenodo.org/records/10256836.
- Gemma Team, Thomas Mesnard, Cassidy Hardin, Robert Dadashi, Surya Bhupatiraju, Shreya Pathak, Laurent Sifre, Morgane Rivière, Mihir Sanjay Kale, Juliette Love, Pouya Tafti, Léonard Hussenot, Pier Giuseppe Sessa, Aakanksha Chowdhery, Adam Roberts, Aditya Barua, Alex Botev, Alex Castro-Ros, Ambrose Slone, Amélie Héliou, Andrea Tacchetti, Anna Bulanova, Antonia Paterson, Beth Tsai, Bobak Shahriari, Charline Le Lan, Christopher A. Choquette-Choo, Clément Crepy, Daniel Cer, Daphne Ippolito, David Reid, Elena Buchatskaya, Eric Ni, Eric Noland, Geng Yan, George Tucker, George-Christian Muraru, Grigory Rozhdestvenskiy, Henryk Michalewski, Ian Tenney, Ivan Grishchenko, Jacob Austin, James Keeling, Jane Labanowski, Jean-Baptiste Lespiau, Jeff Stanway, Jenny Brennan, Jeremy Chen, Johan Ferret, Justin Chiu, Justin Mao-Jones, Katherine

540 Lee, Kathy Yu, Katie Millican, Lars Lowe Sjoesund, Lisa Lee, Lucas Dixon, Machel Reid, Maciej 541 Mikuła, Mateo Wirth, Michael Sharman, Nikolai Chinaev, Nithum Thain, Olivier Bachem, Oscar 542 Chang, Oscar Wahltinez, Paige Bailey, Paul Michel, Petko Yotov, Rahma Chaabouni, Ramona 543 Comanescu, Reena Jana, Rohan Anil, Ross McIlroy, Ruibo Liu, Ryan Mullins, Samuel L Smith, 544 Sebastian Borgeaud, Sertan Girgin, Sholto Douglas, Shree Pandya, Siamak Shakeri, Soham De, Ted Klimenko, Tom Hennigan, Vlad Feinberg, Wojciech Stokowiec, Yu hui Chen, Zafarali Ahmed, 545 Zhitao Gong, Tris Warkentin, Ludovic Peran, Minh Giang, Clément Farabet, Oriol Vinyals, Jeff 546 Dean, Koray Kavukcuoglu, Demis Hassabis, Zoubin Ghahramani, Douglas Eck, Joelle Barral, 547 Fernando Pereira, Eli Collins, Armand Joulin, Noah Fiedel, Evan Senter, Alek Andreev, and 548 Kathleen Kenealy. Gemma: Open models based on gemini research and technology. arxiv preprint: 549 arXiv:2403.08295, 2024. 550 551 Dan Hendrycks, Collin Burns, Steven Basart, Andy Zou, Mantas Mazeika, Dawn Song, and Jacob 552 Steinhardt. Measuring massive multitask language understanding. In ICLR, 2021. 553 Edward J Hu, yelong shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, 554 and Weizhu Chen. LoRA: Low-rank adaptation of large language models. In ICLR, 2022. 555 Patrick Lewis, Ethan Perez, Aleksandra Piktus, Fabio Petroni, Vladimir Karpukhin, Naman Goyal, 556 Heinrich Küttler, Mike Lewis, Wen-tau Yih, Tim Rocktäschel, et al. Retrieval-augmented generation for knowledge-intensive nlp tasks. In NeurIPS, 2020. 558 559 Quentin Lhoest, Albert Villanova del Moral, Yacine Jernite, Abhishek Thakur, Patrick von Platen, Suraj Patil, Julien Chaumond, Mariama Drame, Julien Plu, Lewis Tunstall, Joe Davison, Mario 561 Šaško, Gunjan Chhablani, Bhavitvya Malik, Simon Brandeis, Teven Le Scao, Victor Sanh, Canwen 562 Xu, Nicolas Patry, Angelina McMillan-Major, Philipp Schmid, Sylvain Gugger, Clément Delangue, 563 Théo Matussière, Lysandre Debut, Stas Bekman, Pierric Cistac, Thibault Goehringer, Victor Mustar, François Lagunas, Alexander Rush, and Thomas Wolf. Datasets: A community library for natural language processing. In EMNLP: System Demonstrations, 2021. 565 566 Xianming Li and Jing Li. Angle-optimized text embeddings. ACL, 2023. 567 Yuanzhi Li, Sébastien Bubeck, Ronen Eldan, Allie Del Giorno, Suriya Gunasekar, and Yin Tat Lee. 568 Textbooks are all you need ii: phi-1.5 technical report. arXiv:2309.05463, 2023. 569 570 Haokun Liu, Derek Tam, Muqeeth Mohammed, Jay Mohta, Tenghao Huang, Mohit Bansal, and Colin 571 Raffel. Few-shot parameter-efficient fine-tuning is better and cheaper than in-context learning. In 572 NeurIPS, 2022. 573 Peter J. Liu, Mohammad Saleh, Etienne Pot, Ben Goodrich, Ryan Sepassi, Lukasz Kaiser, and Noam 574 Shazeer. Generating wikipedia by summarizing long sequences. In *ICLR*, 2018. 575 576 Shih-Yang Liu, Chien-Yi Wang, Hongxu Yin, Pavlo Molchanov, Yu-Chiang Frank Wang, Kwang-577 Ting Cheng, and Min-Hung Chen. Dora: Weight-decomposed low-rank adaptation. In ICML, 2024. 578 579 Sourab Mangrulkar, Sylvain Gugger, Lysandre Debut, Younes Belkada, Sayak Paul, and Benjamin 580 Bossan. PEFT: State-of-the-art parameter-efficient fine-tuning methods. https://github. 581 com/huggingface/peft, 2022. 582 Mistral AI. Codestral: Hello, World! — mistral.ai. https://mistral.ai/news/ 583 codestral/, 2024. [Accessed 01-10-2024]. 584 585 OpenAI. Gpt-4 technical report. arXiv:2303.08774, 2024. 586 Long Ouyang, Jeff Wu, Xu Jiang, Diogo Almeida, Carroll L. Wainwright, Pamela Mishkin, Chong 587 Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, John Schulman, Jacob Hilton, Fraser Kelton, 588 Luke Miller, Maddie Simens, Amanda Askell, Peter Welinder, Paul Christiano, Jan Leike, and 589 Ryan Lowe. Training language models to follow instructions with human feedback. In NeurIPS, 590 2022a. Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong 592 Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, et al. Training language models to follow instructions with human feedback. In NeurIPS, 2022b.

594 595 596	Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J. Liu. Exploring the limits of transfer learning with a unified text-to-text transformer. <i>JMLR</i> , 2020.
597 598 599 600	Baptiste Roziere, Jonas Gehring, Fabian Gloeckle, Sten Sootla, Itai Gat, Xiaoqing Ellen Tan, Yossi Adi, Jingyu Liu, Tal Remez, Jérémy Rapin, et al. Code llama: Open foundation models for code. <i>arXiv:2308.12950</i> , 2023.
601 602	Maarten Sap, Hannah Rashkin, Derek Chen, Ronan Le Bras, and Yejin Choi. Social iqa: Common- sense reasoning about social interactions. In <i>EMNLP</i> , 2019.
603 604	Mike Schuster and Kuldip K. Paliwal. Bidirectional recurrent neural networks. In <i>IEEE Transactions</i> on Signal Processing, 1997.
605 606 607	Jacob Mitchell Springer, Suhas Kotha, Daniel Fried, Graham Neubig, and Aditi Raghunathan. Repetition improves language model embeddings. <i>arviv preprint arXiv2402.15449</i> , 2024.
608 609 610	Alon Talmor, Jonathan Herzig, Nicholas Lourie, and Jonathan Berant. CommonsenseQA: A question answering challenge targeting commonsense knowledge. In Jill Burstein, Christy Doran, and Thamar Solorio (eds.), <i>NAACL</i> , 2019.
611 612 613	Rohan Taori, Ishaan Gulrajani, Tianyi Zhang, Yann Dubois, Xuechen Li, Carlos Guestrin, Percy Liang, and Tatsunori B. Hashimoto. Stanford alpaca: An instruction-following llama model. https://github.com/tatsu-lab/stanford_alpaca, 2023.
614 615 616 617	Yi Tay, Mostafa Dehghani, Vinh Q. Tran, Xavier Garcia, Jason Wei, Xuezhi Wang, Hyung Won Chung, Dara Bahri, Tal Schuster, Steven Zheng, Denny Zhou, Neil Houlsby, and Donald Metzler. UL2: Unifying language learning paradigms. In <i>ICLR</i> , 2023.
618 619	Arun James Thirunavukarasu, Darren Shu Jeng Ting, Kabilan Elangovan, Laura Gutierrez, Ting Fang Tan, and Daniel Shu Wei Ting. Large language models in medicine. <i>Nature medicine</i> , 2023.
620 621 622 623 624 625 626 627 628 629 630 631 632	Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, Dan Bikel, Lukas Blecher, Cristian Canton Ferrer, Moya Chen, Guillem Cucurull, David Esiobu, Jude Fernandes, Jeremy Fu, Wenyin Fu, Brian Fuller, Cynthia Gao, Vedanuj Goswami, Naman Goyal, Anthony Hartshorn, Saghar Hosseini, Rui Hou, Hakan Inan, Marcin Kardas, Viktor Kerkez, Madian Khabsa, Isabel Kloumann, Artem Korenev, Punit Singh Koura, Marie-Anne Lachaux, Thibaut Lavril, Jenya Lee, Diana Liskovich, Yinghai Lu, Yuning Mao, Xavier Martinet, Todor Mihaylov, Pushkar Mishra, Igor Molybog, Yixin Nie, Andrew Poulton, Jeremy Reizenstein, Rashi Rungta, Kalyan Saladi, Alan Schelten, Ruan Silva, Eric Michael Smith, Ranjan Subramanian, Xiaoqing Ellen Tan, Binh Tang, Ross Taylor, Adina Williams, Jian Xiang Kuan, Puxin Xu, Zheng Yan, Iliyan Zarov, Yuchen Zhang, Angela Fan, Melanie Kambadur, Sharan Narang, Aurelien Rodriguez, Robert Stojnic, Sergey Edunov, and Thomas Scialom. Llama 2: Open foundation and fine-tuned chat models. <i>arXiv:2307.09288</i> , 2023.
633 634	Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Ł ukasz Kaiser, and Illia Polosukhin. Attention is all you need. In <i>NeurIPS</i> , 2017.
635 636 637	Thomas Wang, Adam Roberts, Teven Le Scao Daniel Hesslow, Hyung Won Chung, Iz Beltagy, Julien Launay, and Colin Raffel. What language model architecture and pretraining objective work best for zero-shot generalization? In <i>ICML</i> , 2022.
638 639 640 641 642	Thomas Wolf, Lysandre Debut, Victor Sanh, Julien Chaumond, Clement Delangue, Anthony Moi, Pierric Cistac, Tim Rault, Rémi Louf, Morgan Funtowicz, Joe Davison, Sam Shleifer, Patrick von Platen, Clara Ma, Yacine Jernite, Julien Plu, Canwen Xu, Teven Le Scao, Sylvain Gugger, Mariama Drame, Quentin Lhoest, and Alexander M. Rush. Transformers: State-of-the-art natural language processing. In <i>EMNLP: System Demonstrations</i> , 2020.
643 644 645	Zhilin Yang, Zihang Dai, Yiming Yang, Jaime Carbonell, Russ R Salakhutdinov, and Quoc V Le. Xlnet: Generalized autoregressive pretraining for language understanding. <i>NeurIPS</i> , 2019.
646 647	Shengyu Zhang, Linfeng Dong, Xiaoya Li, Sen Zhang, Xiaofei Sun, Shuhe Wang, Jiwei Li, Runyi Hu, Tianwei Zhang, Fei Wu, and Guoyin Wang. Instruction tuning for large language models: A survey. <i>arXiv:2308.10792</i> , 2024.

```
648
       6
           APPENDIX
649
650
651
       6.1 PSEUDOCODE FOR BITUNE TRAINING STEP
652
653
654
655
656
       Algorithm 2 Python-like pseudocode of Bitune Training Step.
657
       #
         prompt - tensor with tokenized instruction
658
         answer - tensor with tokenizer answer to model
       #
659
       # theta - trainable tensor with mixing coeff. for each layer
660
       # theta_init - initial value of mixing coefficients
661
       k_c, v_c = model_causal(prompt) # Pass on the instruction for causal
662
           features
663
       k_b, v_b = model_bidir(prompt) # Pass to obtain bidirectional features
       # Combine both sets of features
664
       alpha = theta.abs() / (theta.abs() + theta_init)
665
       k = k_c * (1 - alpha) + k_b * alpha
666
       v = v_c * (1 - alpha) + v_b * alpha
667
       kv = (k, v)
668
       logits = model_causal(answer, kv)
669
       loss = compute_loss(logits, answer)
670
       loss.backward()
       update_parameters(model_causal, model_bidir, theta)
671
672
673
674
675
       6.2 HYPERPARAMETERS
676
677
678
679
680
                       Table 7: Hyperparameters shared across models and datasets.
681
682
                          Hyperparameter
                                                              Value
683
                          GPUs
                                                                1
684
                                                             AdamW
                          Optimizer
685
                          LR Scheduler
                                                             Linear
686
                          Weight Decay
                                                               0.0
687
                          Batch Size (incl. accumulation)
                                                               10
                                                               10
                          Accumulation Steps
688
                                                       10% of update steps
                          Warmup Steps
689
                          \theta_{\text{init}} (Bitune)
                                                              0.01
690
                          RNG Seeds
                                                            42, 43, 44
691
692
693
694
695
696
697
                               Table 8: Dataset-specific hyperparameters.
699
               Hyperparameter
                               UltraFeedback
                                              PIQA
                                                      ARC
                                                             CSQA
                                                                     SIQA
                                                                            GSM8K
700
               Epochs
                                     3
                                                1
                                                       5
                                                               1
                                                                       1
                                                                                1
701
                                   3000
                                               1605
                                                              974
                                                                              747
                                                      555
                                                                     3341
               Update Steps
```

Model	UltraFeedback	PIQA	ARC	CSQA	SIQA	GSM8K
Gemma-2B	3E-4	1E-3	1E-3	1E-3	1E-3	1E-3
Gemma-7B	3E-4	3E-4	3E-4	3E-4	3E-4	3E-4
Llama2-7B	3E-4	1E-3	1E-3	1E-3	1E-3	1E-3
Llama3-8B	3E-4	1E-3	1E-3	1E-3	1E-3	1E-3
Phi-2	3E-4	1E-3	1E-3	1E-3	1E-3	1E-3
Gemma-2B (DoRA)	3E-4	-	-	-	-	-
Gemma-2B (IA3)	1E-3	-	-	-	-	-
Llama3-8B (DoRA)	3E-4	-	-	-	-	-
Llama3-8B (IA3)	1E-3	-	-	-	-	-

Table 9: Learning rate for given dataset-model pair, including different PEFT variants for instruction tuning experiments.

Table 10: Configuration of PEFT methods. All other hyperparameters have default values of HuggingFace PEFT library.

Hyperparameter	Value
Rank (LoRA, DoRA)	8
Alpha (LoRA, DoRA)	1
Target Modules (All)	All linear layers of MLP and Self-Attention
Feedforward Modules (IA3)	All linear layers of MLP

#### 6.3 DATASETS

For all experiments we used HuggingFace Datasets (Lhoest et al., 2021) library to obtain necessary datasets.

Table 11: Table with datasets and corresponding *paths*, to be used with HuggingFace Datasets library.

Dataset	Path
UltraFeedback	openbmb/UltraFeedback
PIQA	piqa
ARC-Challenge	allenai/ai2_arc
CSQA	tau/commonsense_qa
SIQA	social_i_qa
GSM8K	gsm8k
MMLU	hails/mmlu_no_train

742 6.4 Training Speed & Memory Usage

As the method introduces two additional forward passes during training, both the training speed and the memory usage are impacted. Here we present average training times and GPU memory usage on the instruction-tuning setup with 3000 update steps (30000 actual steps, due to gradient accumulation), on a single A100 GPU, for models of two different scales - Gemma-2B & Llama3-8B. Our implementation has not been optimized, which means that e.g. training times could be improved via parallelization of two passes on the prompt. Table shows average training time, peak GPU memory usage during training, and average accuracy on downstream tasks. Additionally, we provide inference times, averaged over 10 runs, for a given prompt-to-answer length in tokens. 

51							
50			Train. Time [h]	Memory [GB]	Acc.	Inference	e Time [s]
753	Model	Method		•		50:200	200:50
754	Gemma-2B	LoRA	1.0	14.9	45.7	6.31	1.56
755		Bitune	3.1	19.8	49.7	6.72	1.74

756	Llama3-8B	LoRA	1.7	26.6	70.4	11.48	2.39
757		Bitune	5.3	30.8	72.7	11.65	2.73
750		Dituite	5.5	50.0	12.1	11.05	2.15

Using these values, one can approximate required compute to reproduce results on a given tasks, as all experiments shared the same batch size and many other hyperparameters.

#### 6.5 INSTRUCTION-TUNING ON ALPACA DATASET

We tested Bitune on another, larger instruction-tuning dataset - cleaned Alpaca dataset<sup>1</sup> (Taori et al., 2023) with 50,000 samples. Similarly to previous experimental settings, first we tested different learning rates for the baseline LoRA finetuning, picked the best one, and then used the same learning rate for other approaches used in the experiment - Bitune, and Naive Bidir. (introduced in the ablation study section 3.3). The results demonstrate that our method's benefits extend to larger datasets as well.

Table 13: Zero-shot results after instruction-tuning on the Alpaca dataset. Results are averaged over 3 random seeds.

		PIQA	ARC	CSQA	SIQA	MMLU	Avg.
Model	Method						
Gemma-2B	LoRA	64.1	40.8	38.6	38.6	33.8	43.2
	Naive Bidir.	64.9	41.7	42.2	42.3	33.6	44.9
	Bitune	66.5	43.6	43.2	43.5	36.2	46.6
Llama3-8B	LoRA	76.6	71.1	63.3	67.4	57.8	67.2
	Naive Bidir.	79.5	70.5	62.4	66.8	55.3	66.9
	Bitune	79.4	71.3	63.5	68.6	58.3	68.2

#### 6.6 FINETUNING LARGER MODEL ON GSM8K

In order to verify whether improvements hold for larger, already highly capable models, we finetune
Codestral (Mistral AI, 2024) with 22B parameters on the GSM8K dataset. The results, averaged
over 3 random seeds, show a substantial 4.3 percentage point improvement over the baseline LoRA
finetuning, indicating that Bitune can provide considerable gains even for larger models, which
already have strong performance on a given task.

Table 14: Results for Bitune and baseline LoRA finetuning on GSM8K dataset.

		GSM8K
Model	Method	
Codestral-22B	LoRA	69.3
	Bitune	73.6

#### 6.7 PROMPT TEMPLATES

Templates used to format instruction-answer pairs for a given dataset, for both training and evaluation. In all cases there is a space character at the beginning of the answer part.

Dataset	Instruction	Answer
UltraFeedback	Question: {instruction}	{completion} <eos></eos>
	Answer:	

<sup>1</sup>https://huggingface.co/datasets/yahma/alpaca-cleaned

PIQA	Question: {question}	{answer} <eos></eos>
	Choices:	
	{choice0}	
	{choice1}	
	Answer:	
ARC	Question: {question}	{answer} <eos></eos>
	Choices:	
	{choice0}	
	{choice1}	
	{choice2}	
	{choice3}	
	Answer <sup>.</sup>	
CSQA	Question: {question}	{answer} <eos></eos>
	Choices:	
	{cnoice0}	
	{cnoice1}	
	{cnoice2}	
	{choice5}	
	{choice4}	
	Answer:	
SIQA	Question: Given the context, answer correctly the	({answer_index}) <eos></eos>
	question.	
	Context: {context}	
	Question: {question}	
	Chairman	
	(0) (1) (1) (1) (1) (1) (1) (1) (1) (1) (1	
	$(0) \{\text{cnoice0}\}$	
	$(1) \{ \text{choice} 1 \}$	
	(2) {choice2}	
	Answer:	
GSM8K	Question: {question}	{answer} <eos></eos>
	Answer:	
MMLU	Question: {question}	({answer_index}) <eos></eos>
	(0). {choice0}	
	(1). {choice1}	
	(2). {choice2}	
	(3). {choice3}	
	Answer:	
0 <b>D</b>		
.8 RESULTS	S WITH STANDARD DEVIATION	
ables with cor	nplete results averaged over 3 seeds, includes standard of	deviation.
	1 U I I I I I I I I I I I I I I I I I I	

							0				
			PIQA		ARC	(	CSQA		SIQA	Μ	MI
		mean	std	mean	std	mean	std	mean	std	mean	st
Model	Method										
Gemma-2B	Pretrained	57.51	-	36.86	-	35.46	-	38.18	-	34.01	-
	LoRA	66.74	0.52	43.37	1.33	42.32	1.60	44.27	3.85	31.74	0
	LoRA <sub>16</sub>	66.50	0.85	42.72	1.46	42.34	1.21	43.82	3.71	31.61	0
	Bitune	69.59	1.20	47.47	0.64	46.87	1.98	49.51	1.54	35.29	0
Gemma-7B	Pretrained	73.12	-	78.33	-	62.00	-	64.74	-	59.04	-
	LoRA	84.24	0.51	79.18	0.76	68.47	1.61	71.92	0.34	55.26	0
	LoRA <sub>16</sub>	83.93	0.94	79.24	0.61	68.39	2.05	71.99	0.83	53.41	4
	Bitune	83.59	0.46	80.09	0.90	69.15	0.56	72.74	0.90	53.81	0
Llama2-7B	Pretrained	59.25	-	38.14	-	32.60	-	45.09	-	35.98	-
	LoRA	69.51	0.73	49.94	0.79	45.32	2.65	57.05	1.24	41.06	0
	LoRA <sub>16</sub>	69.86	0.25	49.89	1.51	45.59	2.08	56.69	1.89	41.21	0
	Bitune	70.00	0.53	51.11	0.23	48.08	2.59	59.09	0.96	41.87	C
Llama3-8B	Pretrained	68.99	-	73.63	-	65.36	-	56.81	-	56.00	-
	LoRA	81.94	0.38	74.46	1.06	69.23	0.67	68.99	1.55	57.62	0
	LoRA <sub>16</sub>	82.35	0.83	74.91	0.45	70.52	0.46	68.63	1.38	57.98	0
	Bitune	84.39	0.24	77.42	1.15	72.70	0.82	70.15	0.34	58.96	C
Phi-2	Pretrained	70.35	-	67.32	-	61.43	-	65.05	-	45.41	-
	LoRA	76.31	0.17	66.67	0.44	61.62	0.40	66.65	0.31	48.23	0
	LoRA <sub>16</sub>	76.12	0.93	66.61	0.27	61.64	0.45	66.82	0.83	47.73	C
	Bitune	76.51	0.30	67.18	0.49	63.04	0.71	68.51	0.56	48.92	0
				1		1		1		1	

Table 16: Zero-shot results on downstream tasks after instruction-tuning on the UltraFeedback dataset.

Table 17: Result for the downstream task training setup.

005				DIO		4 D.C		70.01		0101		
905				PIQA		ARC	(	LSQA		SIQA	GS	SM8K
906			mean	std	mean	std	mean	std	mean	std	mean	std
907	Model	Method										
908	Gemma-2B	LoRA	81.41	0.46	58.05	1.16	77.18	0.05	77.40	0.40	30.17	0.75
909		Bitune	83.28	0.33	59.98	1.19	78.32	0.68	76.60	0.24	32.98	0.97
910	Phi-2	LoRA	82.79	1.15	76.31	0.56	78.68	0.76	80.30	0.44	58.55	0.89
911		Bitune	83.91	0.28	77.02	0.34	78.95	0.22	80.38	0.21	59.21	1.15
912	Gemma-7B	LoRA	91.42	0.60	84.58	0.36	84.41	1.04	79.36	1.41	59.11	0.97
913		Bitune	92.13	0.27	84.22	0.74	84.22	0.76	79.39	1.33	59.39	0.72
914	Llama3-8B	LoRA	90.24	0.27	80.75	0.36	83.92	0.74	83.15	0.41	60.45	0.72
915		Bitune	90.52	0.78	81.26	0.79	84.11	1.05	82.12	0.30	63.43	0.27
916	Llama2-7B	LoRA	84.39	0.69	66.55	0.78	81.52	0.83	82.67	0.31	32.02	0.68
917		Bitune	84.39	0.61	66.87	0.90	81.95	0.56	81.39	0.15	32.85	1.05

			PIQA		ARC		CSQA		SIQ	A	Μ
		mean	n std	mean	std	mean	std	mean	sto	d m	ean
Model	Method										
Gemma-2B	LoRA	66.74	4 0.52	43.37	1.33	42.32	1.60	44.27	3.8	35 31	.74
	Bitune	69.59	9 1.20	47.47	0.64	46.87	1.98	49.51	1.5	54 35	5.29
	Naive Bidir.	67.79	9 1.00	44.65	1.78	46.79	2.86	48.04	0.9	95 32	2.43
	No Mixing	69.01	1.47	45.71	1.32	46.14	2.56	49.80	0.5	53 34	.03
	Only Causal	66.39	9 1.28	45.28	0.65	42.45	1.65	46.98	0.7	6 33	5.27
	Shared Weights	68.10	0.41	44.34	0.90	44.53	1.67	47.19	2.1	.6 32	2.93
Llama3-8B	LoRA	81.94	4 0.38	74.46	1.06	69.23	0.67	68.99	1.5	55 57	.62
	Bitune	84.39	0.24	77.42	1.15	72.70	0.82	70.15	0.3	<sup>34</sup> 58	8.96
	Naive Bidir.	85.44	4 0.25	76.45	0.60	69.37	0.28	70.04	0.4	3 58	3.33
	No Mixing	85.56	6 0.38	74.86	0.21	68.60	1.31	69.74	0.1	1 58	8.88
	Only Causal	82.37	0.41	74.63	0.91	70.65	0.93	69.17	1.5	51 58	3.43
	Shared weights	84.10	0.63	/5.94	1.11	/1.91	0.71	/0.61	0.6	0 59	0.03
1401	e 19: Results for d	ifferent PI	PEFT m	ethods u	sed in c $\overline{C}$	ombinat	ion with	n Bitune	DA	N	4MI
	e 19: Results for d	ifferent PI ean	PEFT m QA std m	ethods u AR ean st	sed in c	combinat CSQ ean st	ion with A d me	n Bitune SIC	A QA td	Mmean	1ML st
Model	m Method	ifferent PI ean	PEFT m QA std m	ethods u AR ean st	sed in c C d mo	combinat CSQ ean st	ion with A d me	n Bitune SIC ean s	A QA td	M mean	1ML st
Model Gemma-2B	m Method LoRA 66	PI ean	QA std m	AR AR ean st	sed in c RC rd ma 33 42	CSQ ean st	ion with A d me	SIC SIC ean s	2A td .85	N mean 31.74	4ML st
Model Gemma-2B	m Method LoRA 66 Bitune 69	Ifferent           PI           ean           5.74           0.59	QA std m 0.52 43 1.20 47	AR ean st 3.37 1.3 7.47 0.0	sed in c RC rd mo 33 42 64 46	CSQ ean st 	ion with A d me 50 44 98 49	SIC ean s .27 3. .51 1.	2A td .85 .54	N mean 31.74 35.29	4ML st
Model Gemma-2B	m Method LoRA 66 Bitune 69 DoRA 66	ifferent           PI           ean           5.74           0.59           1           5.70	PEFT m QA std m 0.52 43 1.20 47 0.63 43	AR ean st 3.37 1 7.47 0.0 9.57 0.9	sed in c RC ad mo 33 42 64 46 91 41	CSQ ean st .32 1.0 .87 1.9 .88 1.0	ion with DA d me 50 44 98 49 08 44	Bitune           SIC           ean         s           .27         3.           .51         1.           .71         3.	QA td .85 .54 .83	N mean 31.74 35.29 31.95	4ML st 0.7 0.0 0.7
Model Gemma-2B	m Method LoRA 66 Bitune 69 DoRA 66 Bitune <sub>DoRA</sub> 69	PI           ean           5.74           0.59           1           5.70           0.62	QA std m 0.52 43 1.20 47 0.63 43 0.69 47	AR ean st 3.37 1.1 7.47 0.0 3.57 0.9 7.47 0.1	sed in c C 33 42 64 46 91 41 34 46	CSQ ean st 	ion with A ma 50 44 98 49 08 44 42 49	Bitune           SIC           ean         s           .27         3.           .51         1.           .71         3.	QA td .85 .54 .83 .80	N mean 31.74 35.29 31.95 35.09	1ML st 0.7 0.0 0.7 0.1
Model Gemma-2B	m Method LoRA 66 Bitune 69 DoRA 66 Bitune <sub>DoRA</sub> 69 IA3 67	PI           ean           5.74           0.59           5.70           0.62           0.62	PEFT m           QA           std         m           0.52         43           1.20         47           0.63         43           0.69         47           0.52         46	AR           aean         st           3.37         1.3           3.37         1.4           3.37         0.4           3.57         0.9           3.57         0.9           3.57         0.9           5.57         0.9           5.50         0.0	sed in c C 33 42 64 46 91 41 34 46 67 45	CSQ ean st .32 1.0 .87 1.9 .88 1.0 .87 2.4 .54 0.0	ion with A ma 50 44 98 49 08 44 42 49 30 37	SIC           can         s           .27         3.           .51         1.           .71         3.           .71         1.           .65         1.	2A td .85 .54 .83 .80 .07	W mean 31.74 35.29 31.95 35.09 32.49	1ML st 0.7 0.0 0.7 0.1 0.2
Model Gemma-2B	m Method LoRA 66 Bitune 69 DoRA 66 Bitune <sub>DoRA</sub> 69 IA3 67 Bitune <sub>IA3</sub> 67	PI           ean           5.74         0           5.74         0           5.74         0           5.70         0           0.62         0           7.25         0           7.54         1	PEFT m           QA           std         m           0.52         43           1.20         47           0.63         43           0.52         46           1.27         47	AR           action         AR           action         str           3.37         1.3           7.47         0.6           3.57         0.9           7.47         0.1           5.57         0.9           7.47         0.1           5.50         0.6           7.27         0.3	sed in c C 33 42 64 46 91 41 34 46 67 45 53 48	combinat CSQ ean st .32 1.0 .87 1.9 .88 1.0 .87 2.4 .54 0.1 .92 0.9	ion with A d me 50 44 98 49 98 44 42 49 30 37 94 44	SIC           ean         s           .27         3.           .51         1.           .71         3.           .71         1.           .65         1.           .29         1.	A. QA td .85 .54 .83 .80 .07 .08	N mean 31.74 35.29 31.95 35.09 32.49 33.60	1ML sto 0.7 0.0 0.7 0.1 0.2 0.3
Model Gemma-2B Llama3-8B	m Method LoRA 66 Bitune 69 DoRA 66 Bitune <sub>DoRA</sub> 69 IA3 67 Bitune <sub>IA3</sub> 67 LoRA 81	PI           ean           5.74         0           5.9         1           5.70         0           0.62         0           7.25         0           7.54         1          94         0	PEFT m         QA         std       m         0.52       43         1.20       47         0.63       43         0.69       47         0.52       46         1.27       47         0.38       74	AR           aean         st           3.37         1.3           7.47         0.4           3.57         0.9           7.47         0.1           5.50         0.6           7.27         0.3           4.46         1.6	sed in c C 33 42 64 46 91 41 34 46 67 45 53 48 06 69	CSQ           ean         st           .32         1.0           .87         1.9           .88         1.0           .87         2.4           .92         0.9           .23         0.0	ion with A d me 60 44 98 49 08 44 42 49 30 37 94 44 57 68	SIC           .27         3.           .51         1.           .71         3.           .71         1.           .65         1.           .29         1.           .99         1.	2A td .54 .83 .80 .07 .08 .55	N mean 31.74 35.29 31.95 35.09 32.49 33.60 57.62	4ML st 0.7 0.0 0.1 0.2 0.3 0.5
Model Gemma-2B Llama3-8B	m Method LoRA 66 Bitune 69 DoRA 66 Bitune <sub>DoRA</sub> 69 IA3 67 Bitune <sub>IA3</sub> 67 LoRA 81 Bitune 84	PI           ean           5.74         0           0.59         1           5.70         0           0.62         0           7.25         0           7.54         1           .94         0           .39         0	PEFT m           QA           std         m           0.52         43           1.20         47           0.63         43           0.69         47           0.52         46           1.27         47           0.38         74           0.24         77	AR           aean         st           3.37         1           2.47         0.0           3.57         0.9           2.47         0.0           3.50         0.4           5.50         0.4           4.46         1.0           4.42         1.1	sed in c C 33 42 64 46 91 41 34 46 67 45 53 48 06 69 15 72	CSQ           ean         st           .32         1.0           .87         1.9           .88         1.0           .87         2.4           .54         0.1           .92         0.9           .23         0.0           .70         0.3	ion with A d me 50 44 98 49 908 44 42 49 30 37 94 44 57 68 82 70	SIC           ean         s           .27         3.           .51         1.           .71         3.           .71         1.           .65         1.           .99         1.           .15         0.	2A td 	N mean 31.74 35.29 31.95 35.09 32.49 33.60 57.62 58.96	IML         st         0.1         0.2         0.3         0.4         0.5         0.5
Model Gemma-2B Llama3-8B	m Method LoRA 66 Bitune 69 DoRA 66 Bitune <sub>DoRA</sub> 69 IA3 67 Bitune <sub>IA3</sub> 67 LoRA 81 Bitune 84 DoRA 82	PI           ean           5.74         0           0.59         1           5.70         0           0.62         0           7.54         1          94         0           1.39         0	PEFT m           QA           std         m           0.52         43           1.20         47           0.63         43           0.69         47           0.52         46           1.27         47           0.38         74           0.24         77           0.44         75	AR           aan         st           3.37         1           3.37         1           3.37         1           3.37         1           3.37         0           3.57         0.9           3.57         0.9           3.57         0.9           3.57         0.9           3.50         0.0           3.50         0.0           3.47         0.3           4.46         1.0           4.46         1.0           4.42         1.           5.40         0.3	sed in c C 33 42 64 46 91 41 34 46 67 45 53 48 06 69 15 72 52 70	CSQ           ean         st           .32         1.0           .87         1.9           .88         1.0           .87         2.4           .54         0.3           .92         0.9           .23         0.0           .70         0.3           .22         0.3	ion with A d me 50 44 98 49 08 44 42 49 30 37 94 44 57 68 82 70 39 69	SIC           ean         s           .27         3.           .51         1.           .71         3.           .71         1.           .65         1.           .99         1.           .15         0.           .19         1.	QA           td           .85           .54           .83           .80           .07           .08           .55           .34           .14	W mean 31.74 35.29 31.95 35.09 32.49 33.60 57.62 58.96 57.67	1ML           str           0.7           0.0           0.7           0.1           0.2           0.3           0.4           0.5           0.2           0.3
Model Gemma-2B Llama3-8B	m Method LoRA 66 Bitune 69 DoRA 66 Bitune <sub>DoRA</sub> 69 IA3 67 Bitune <sub>IA3</sub> 67 LoRA 81 Bitune 84 DoRA 82 Bitune <sub>DoRA</sub> 84	PI           ean           5.74         0           0.59         1           5.70         0           0.62         0           7.25         0           7.54         1	PEFT m           QA           std         m           0.52         43           1.20         47           0.63         43           0.69         47           0.52         46           1.27         47           0.38         74           0.24         77           0.44         75           0.46         77	AR           ean         st           3.37         1           7.47         0.0           3.57         0.9           7.47         0.0           3.57         0.9           7.47         0.1           5.50         0.0           7.27         0.3           4.46         1.0           7.42         1.3           6.40         0.3	sed in c C 33 42 64 46 91 41 34 46 67 45 53 48 06 69 15 72 52 70 37 71 1	CSQ           ean         st           .32         1.0           .87         1.9           .88         1.0           .87         2.4           .54         0.3           .92         0.1           .23         0.0           .22         0.1           .99         0.4	ion with A d me 50 44 98 49 08 44 42 49 30 37 94 44 57 68 82 70 39 69 57 70	A Bitune           SIC           ean         s           .27         3.           .51         1.           .71         3.           .71         1.           .65         1.           .99         1.           .15         0.           .19         1.           .56         0.	QA td 	W mean 31.74 35.29 31.95 35.09 32.49 33.60 57.62 58.96 57.67 58.72	IML           str           0.7           0.0           0.7           0.1           0.2           0.3           0.4           0.5           0.2           0.5           0.5
Model Gemma-2B Llama3-8B	m Method LoRA 66 Bitune 69 DoRA 66 Bitune <sub>DoRA</sub> 69 IA3 67 Bitune <sub>IA3</sub> 67 LoRA 81 Bitune 84 DoRA 82 Bitune <sub>DoRA</sub> 84 IA3 80	PI           ean           5.74         0           5.74         0           5.74         0           5.70         0           0.62         0           7.25         0           7.54         1	PEFT m         QA         std       m         0.52       43         1.20       47         0.63       43         0.69       47         0.52       46         1.27       47         0.38       74         0.44       75         0.46       77         0.58       75	AR           ean         st           3.37         1           7.47         0.0           3.57         0.9           7.47         0.0           3.57         0.9           7.47         0.0           5.50         0.0           7.27         0.1           4.46         1.0           7.42         1.1           5.40         0.1           5.54         0.1	sed in c C 33 42 64 46 91 41 34 46 67 45 53 48 06 69 15 72 52 70 37 71 10 68	CSQ           ean         st           .32         1.0           .87         1.9           .88         1.0           .87         2.4           .54         0.3           .92         0.3           .23         0.0           .22         0.3           .99         0.4           .25         0.3	ion with 0 44 0 44 0 44 0 44 0 44 0 44 0 44 1 42 1 49 0 30 1 44 1 44 1 67 1 68 82 1 70 0 39 1 69 1 70 2 6 66 1 1 1 1 1 1 1 1 1 1 1 1 1 1	SIC           can         s           .27         3.           .51         1.           .71         3.           .71         1.           .65         1.           .99         1.           .15         0.           .38         0.	QA td .85 .54 .83 .80 .07 .08 .55 .34 .14 .21 .18	W mean 31.74 35.29 31.95 35.09 32.49 33.60 57.62 58.96 57.67 58.72 58.75	AML           str           0.7           0.1           0.2           0.3           0.4           0.5           0.5           0.5           0.5
Model Gemma-2B Llama3-8B	m Method LoRA 66 Bitune 69 DoRA 66 Bitune <sub>DoRA</sub> 69 IA3 67 Bitune <sub>IA3</sub> 67 LoRA 81 Bitune 84 DoRA 82 Bitune <sub>DoRA</sub> 84 IA3 80 Bitune <sub>IA3</sub> 83	PI           ean           5.74         0           5.74         0           5.74         0           5.70         0           0.62         0           7.25         0           7.54         1           .94         0           8.39         0           0.87         0           8.42         0	PEFT m         QA         std       m         0.52       43         1.20       47         0.63       43         0.69       47         0.52       46         1.27       47         0.38       74         0.44       75         0.46       77         0.58       75         0.52       75	AR           aan         st           3.37         1.3           7.47         0.6           3.57         0.9           7.47         0.6           3.57         0.9           7.47         0.6           5.50         0.6           7.27         0.3           4.46         1.6           7.42         1.3           6.40         0.3           5.54         0.6           6.58         0.4	sed in c C d ma 33 42 64 46 91 41 34 46 67 45 53 48 06 69 15 72 52 70 37 71 10 68 68 69	CSQ           ean         st           .32         1.4           .87         1.9           .88         1.0           .87         2.4           .54         0.4           .92         0.9           .23         0.4           .70         0.5           .92         0.4           .92         0.5           .15         0.4	ion with A d ma 50 44 98 49 98 49 908 44 42 49 30 37 94 44 57 68 82 70 39 69 57 70 26 66 40 67	SIC           ean         s           .27         3.           .51         1.           .71         3.           .71         1.           .65         1.           .99         1.           .15         0.           .38         0.           .83         0.	2A td .85 .54 .83 .80 .07 .08 .55 .34 .14 .21 .18 .25	W mean 31.74 35.29 31.95 35.09 32.49 33.60 57.62 58.96 57.67 58.72 58.75 58.80	AML           st           0.0           0.1           0.2           0.3           0.4           0.5           0.5           0.5           0.5           0.5           0.5           0.5           0.5           0.5           0.5           0.5           0.5           0.5           0.5           0.5           0.5           0.5           0.5
Model Gemma-2B Llama3-8B	m Method LoRA 66 Bitune 69 DoRA 66 Bitune <sub>DoRA</sub> 69 IA3 67 Bitune <sub>IA3</sub> 67 LoRA 81 Bitune 84 DoRA 82 Bitune <sub>DoRA</sub> 84 IA3 80 Bitune <sub>IA3</sub> 83	PI           ean           5.74         0           5.74         0           5.70         0           7.25         0           7.54         1           .94         0           8.39         0           0.87         0           8.42         0	PEFT m         QA         std       m         0.52       43         1.20       47         0.63       43         0.69       47         0.52       46         1.27       47         0.38       74         0.44       75         0.46       77         0.58       75         0.52       75	AR           ean         st           3.37         1.3           2.47         0.4           3.57         0.9           2.47         0.3           2.47         0.3           2.50         0.4           2.27         0.3           2.46         1.0           2.46         1.0           2.46         1.0           2.46         1.0           3.54         0.3           5.54         0.3	sed in c C d ma 33 42 54 46 91 41 34 46 67 45 53 48 06 69 15 72 52 70 37 71 10 68 68 69	CSQ           ean         st           .32         1.4           .87         1.9           .88         1.0           .87         2.4           .54         0.1           .92         0.9           .23         0.0           .22         0.1           .92         0.1           .92         0.1           .92         0.1           .92         0.1           .15         0.4	ion with 24 36 36 37 39 30 37 34 44 30 37 37 34 44 57 68 82 70 39 69 57 70 26 66 40 67 70 26 66 67 70 70 70 70 70 70 70 70 70 7	SIC           can         s           .27         3           .51         1           .71         3           .71         1           .65         1           .29         1           .99         1           .15         0           .38         0           .83         0	QA           td           .85           .54           .83           .80           .07           .08           .55           .34           .14           .21           .18           .25	W mean 31.74 35.29 31.95 35.09 32.49 33.60 57.62 58.96 57.67 58.72 58.75 58.80	IML           st           0.1           0.2           0.3           0.4           0.5           0.6           0.7           0.6           0.7           0.7           0.7           0.7           0.7           0.7           0.7           0.7           0.7           0.7           0.7           0.7
Model Gemma-2B Llama3-8B	m Method LoRA 66 Bitune 69 DoRA 66 Bitune <sub>DoRA</sub> 69 IA3 67 Bitune <sub>IA3</sub> 67 LoRA 81 Bitune 84 DoRA 82 Bitune <sub>DoRA</sub> 84 IA3 80 Bitune <sub>IA3</sub> 83	PI           ean           5.74         0           5.76         1           5.70         0           0.62         0           7.25         0           7.54         1	PEFT m         QA         std       m         0.52       43         1.20       47         0.63       43         0.69       47         0.52       46         1.27       47         0.38       74         0.44       75         0.44       75         0.52       75         0.52       75	AR           ean         st           3.37         1.4           4.47         0.4           5.57         0.9           7.47         0.4           5.50         0.6           7.27         0.3           4.46         1.6           7.43         0.3           5.54         0.5           6.68         0.6	sed in c C d ma 33 42 64 46 91 41 34 46 67 45 53 48 06 69 15 72 52 70 37 71 10 68 68 69	CSQ           ean         st           .32         1.0           .87         1.9           .88         1.0           .87         2.4           .54         0.1           .92         0.9           .23         0.0           .22         0.1           .99         0.1           .15         0.2	ion with A me 50 44 98 49 908 44 42 49 30 37 94 44 57 68 82 70 39 69 57 70 26 66 40 67	SIC           can         s           .27         3           .51         1           .71         3           .71         1           .65         1           .29         1           .99         1           .15         0           .38         0           .83         0	QA           td           .85           .54           .83           .80           .07           .08           .55           .34           .14           .21           .18           .25	M mean 31.74 35.29 31.95 35.09 32.49 33.60 57.62 58.96 57.67 58.72 58.75 58.80	IML           st           0.7           0.6           0.7           0.1           0.2           0.3           0.4           0.5
Model Gemma-2B Llama3-8B	m Method LoRA 66 Bitune 69 DoRA 66 Bitune <sub>DoRA</sub> 69 IA3 67 Bitune <sub>IA3</sub> 67 LoRA 81 Bitune 84 DoRA 82 Bitune <sub>DoRA</sub> 84 IA3 80 Bitune <sub>IA3</sub> 83	PI           ean           0.74         0           0.59         1           5.70         0           0.62         0           7.25         0           7.54         1           .94         0           2.12         0           0.87         0           0.87         0           8.42         0	PEFT m         QA         std       m         0.52       43         1.20       47         0.63       43         0.69       47         0.52       46         1.27       47         0.38       74         0.44       75         0.52       75         0.52       75         0.52       75         ifferent i       1	AR           ean         st           3.37         1.3           7.47         0.4           5.57         0.9           7.47         0.3           5.50         0.6           7.27         0.3           4.46         1.0           4.46         1.0           5.54         0.3           6.54         0.3           6.58         0.6	sed in c C d ma 33 422 64 46 91 41 34 46 67 45 53 48 06 69 15 72 52 70 37 71 10 68 68 69 s for mi	CSQ ean st 	ion with A d ma 60 44 98 49 98 49 908 44 42 49 30 37 94 44 57 68 82 70 39 69 57 70 26 66 40 67 fficients	SIC           can         s           .27         3.           .51         1.           .71         3.           .71         1.           .65         1.           .99         1.           .15         0.           .38         0.           .83         0.	2A td .85 .54 .83 .80 .07 .08 .55 .34 .14 .21 .18 .25	W mean 31.74 35.29 31.95 35.09 32.49 33.60 57.62 58.96 57.67 58.72 58.75 58.80	IML           sta           0.7           0.0           0.7           0.1           0.2           0.3           0.2           0.5           0.2           0.1
Model Gemma-2B Llama3-8B	m Method LoRA 66 Bitune 69 DoRA 66 Bitune <sub>DoRA</sub> 69 IA3 67 Bitune <sub>IA3</sub> 67 LoRA 81 Bitune 84 DoRA 82 Bitune <sub>DoRA</sub> 84 IA3 80 Bitune <sub>IA3</sub> 83 Table 20: Resul	PI           ean           5.74         0           0.59         1           5.70         0           0.62         0           7.25         0           7.54         1	PEFT m         QA         std       m         0.52       43         1.20       47         0.63       43         0.69       47         0.52       46         1.27       47         0.38       74         0.44       75         0.58       75         0.52       75         ifferent i	AR           aan         st           3.37         1           3.37         1           3.37         1           3.37         1           3.37         1           3.37         1           3.37         1           3.37         1           3.37         1           3.37         1           3.57         0.9           2.47         0           5.50         0.0           7.47         0           4.46         1.0           7.42         1           5.40         0           5.54         0           5.68         0.0           nit value:         1	sed in c C d mo 33 42 64 46 91 41 34 46 67 45 53 48 06 69 15 72 52 70 37 71 10 68 68 69 s for mi	combinat CSQ ean st .32 1.0 .87 1.9 .88 1.0 .87 2.4 .54 0.1 .92 0.9 .23 0.0 .70 0.1 .22 0.1 .99 0.0 .15 0.4 xing coe	ion with A d me 50 44 98 49 08 44 42 49 30 37 94 44 57 68 82 70 39 69 57 70 26 66 40 67 fficients	A Bitune           SIC           ean         s           .27         3.           .51         1.           .71         3.           .71         1.           .65         1.           .99         1.           .15         0.           .38         0.           .83         0.	QA         td         .85         .54         .83         .80         .07         .08         .55         .34         .14         .21         .18         .25	W mean 31.74 35.29 31.95 35.09 32.49 33.60 57.62 58.96 57.67 58.72 58.75 58.80	AML           stu           0.7           0.0           0.7           0.1           0.2           0.3           0.2           0.5           0.2           0.1
Model Gemma-2B Llama3-8B	m Method LoRA 66 Bitune 69 DoRA 66 Bitune <sub>DoRA</sub> 69 IA3 67 Bitune <sub>IA3</sub> 67 LoRA 81 Bitune 84 DoRA 82 Bitune <sub>DoRA</sub> 84 IA3 80 Bitune <sub>IA3</sub> 83	PI           ean           5.74         0           5.74         0           5.70         1           5.70         0           0.62         0           7.25         0           7.54         1           .94         0           8.39         0           2.12         0           0.87         0           3.42         0	PEFT m         QA         std       m         0.52       43         1.20       47         0.63       43         0.69       47         0.52       46         1.27       47         0.38       74         0.44       75         0.46       77         0.58       75         0.52       75	AR           aan         st           3.37         1           2.47         0.0           3.57         0.9           2.47         0.1           3.57         0.2           2.47         0.1           3.50         0.4           2.47         0.1           3.50         0.4           4.46         1.0           4.46         1.0           4.46         0.1           4.43         0.1           5.40         0.1           5.54         0.1           5.68         0.0	sed in c 33 42 64 46 91 41 34 46 67 45 53 48 06 69 15 72 52 70 37 71 10 68 68 69	CSQ           ean         st           .32         1.0           .87         1.9           .88         1.0           .87         2.4           .54         0.3           .92         0.9           .23         0.0           .22         0.3           .92         0.4           .92         0.3           .23         0.0           .25         0.3           .15         0.4	ion with A d me 50 44 98 49 908 44 42 49 30 37 94 44 57 68 82 70 39 69 57 70 26 66 40 67	A Bitune           SIC           ean         s           .27         3.           .51         1.           .71         3.           .71         1.           .65         1.           .99         1.           .15         0.           .19         1.           .56         0.           .83         0.	QA           .85           .54           .83           .80           .07           .08           .55           .34           .14           .21           .18           .25	N mean 31.74 35.29 31.95 35.09 32.49 33.60 57.62 58.96 57.67 58.72 58.75 58.80	

			PIQA		ARC	(	LSQA		SIQA	N	MLU
		mean	std								
Model	Init. Value										
Gemma-2B	0.1	69.68	0.74	46.33	0.81	45.97	2.64	49.86	2.21	35.25	0.46
	0.01	69.59	1.20	47.47	0.64	46.87	1.98	49.51	1.54	35.29	0.08
	0.001	67.79	0.88	44.20	0.68	45.56	1.61	46.35	1.08	31.89	0.67
Llama3-8B	0.1	85.00	0.73	76.48	0.63	71.06	1.56	71.68	0.46	59.19	0.38
	0.01	84.39	0.24	77.42	1.15	72.70	0.82	70.15	0.34	58.96	0.37
	0.001	84.60	0.62	76.62	0.09	71.28	1.45	70.62	0.15	58.49	0.22

#### 972 6.9 EXAMPLE ATTENTION MATRICES 973

974 Here we show pairs of matrices with causal and bidirectional attention scores of *Bituned* Gemma-2B
975 for the first two GSM8K samples, with scores averaged over layers and heads. Source tokens are
976 represented with the vertical axis, while target tokens (for which the attention is paid to) with the
977 horizontal one - e.g. the first column represents attention paid to the first (BOS) token, by each other
978 token. Darker color represents higher attention score.



## 1026 6.10 GSM8K SAMPLES

Samples of responses to the first 20 questions from GSM8K benchmark. Shown for Llama3-8B
 intruction-tuned on UltraFeedback dataset - for standard finetuning with LoRA and Bitune.

1030						
1031	Prompt	LoRA	Bitune			
1032	Question: Janet's ducks lay 16 eggs	Janet makes $(3 + 4) = $	Janet makes \$2 per fresh duck egg,			
1033	per day. She eats three for breakfast	day at the farmers' market.	and she sells the remainder of the			
1034	every morning and bakes muffins for		eggs she collects each day. If she eats			
1035	sells the remainder at the farmers'		muffins with four eggs she has 16 - 3			
1036	market daily for \$2 per fresh duck		-4 = 9 eggs left to sell. So, she makes			
1037	egg. How much in dollars does she		$2 \times 9 = 18$ every day at the farmers'			
1038	make every day at the farmers' mar-		market.			
1039	ket?					
1040	Question: A robe takes 2 bolts of	A robe requires 2 bolts of blue fiber	To make a robe, you need 2 bolts of			
1041	blue fiber and half that much white fiber. How many bolts in total does it take?	is equal to 1 bolt of white fiber	fiber So the total number of bolts			
1042		Therefore, the total number of bolts	required is:			
1043		required for a robe is 3 bolts (2 blue	2 bolts of blue fiber + half of 2 bolts			
1044		and 1 white).	of blue fiber (which is 1 bolt of white			
1045			fiber) Tetel helts $-2 + 1 - 2$ helts			
1046			So the total number of bolts required			
1047			to make a robe is 3 bolts.			
1048	Ouestion: Josh decides to try flipping	Josh made a profit of \$50,000.	Josh made a profit of \$20,000.			
1049	a house. He buys a house for \$80,000	Confidence: 95%	To calculate the profit, first find the			
1050	and then puts in \$50,000 in repairs.		increased value of the house after re			
1051	This increased the value of the house		pairs: $\$20,000, 1, \$50,000, -\$120,000$			
1052	make?		580,000 + 550,000 = 5150,000 Then calculate the 150% increase:			
1053	indite :		130,000 * 1.5 = \$195,000			
1054			The difference between the origi-			
1055			nal purchase price and the increased			
1056			value is the profit: \$195,000 - \$80,000 - \$115,000			
1057			Finally, subtract the cost of repairs			
1058			from the profit:			
1059			\$115,000 - \$50,000 = \$20,000			
1060			So, Josh made a profit of \$20,000.			
1061	Question: James decides to run 3	James runs 3 sprints of 60 meters	James runs 3 sprints of 60 meters			
1062	sprints 3 times a week. He runs 60	each, 3 times a week. To find the total	each, 3 times a week. To find the			
1063	meters does he run a week?	ply the number of sprints (3) by the	can multiply the number of sprints			
1064		distance of each sprint (60 meters)	by the distance of each sprint and the			
1065		and then multiply the result by the	number of times he runs the sprints			
1066		number of times he runs the sprints	per week.			
1067		(3) per week. So the total meters lames runs a	10tal meters = $(3 \text{ sprints}) \times (60 \text{ me-})$			
1068		week is:	= 540 meters per week.			
1069		3 sprints x 60 meters per sprint x	So, James runs a total of 540 meters			
1070		3 times per week = 540 meters per	per week.			
1071		week.				
1072						
1073						
1074						
1075						
1076						
1077						

1078

1080	Question: Every day Wendi feeds	To find the number of cups of feed	To find the number of cups of feed
1081	each of her chickens three cups of	Wendi needs to give her chickens in	Wendi needs to give her chickens in
1082	mixed chicken feed, containing seeds,	the final meal of the day, we need	the final meal of the day, we need
1083	mealworms and vegetables to help	to calculate the total amount of feed	to calculate the total amount of feed
1084	keep them healthy. She gives the	needed for the day and divide it by	needed for the day and then divide it
1085	chickens their feed in three separate	the number of meals.	by the number of meals $(3)$ and the
1086	flock of chickens 15 cups of feed. In	the day is:	First let's find the total amount of
1087	the afternoon, she gives her chickens	Morning meal: 15 cups * 20 chickens	feed needed for the day:
1088	another 25 cups of feed. How many	= 300 cups Afternoon meal: 25 cups	Morning meal: 15 cups * 20 chickens
1089	cups of feed does she need to give	* 20 chickens = 500 cups Total feed	= 300 cups Afternoon meal: 25 cups
1090	her chickens in the final meal of the	needed for the day: $300 \text{ cups} + 500$	* 20 chickens = 500 cups Total feed
1091	day if the size of wendi's flock is 20 chickens?	cups = 800 cups Now we divide the total feed needed	needed for the day: $500 \text{ cups} + 500$
1092	entexens.	for the day by the number of meals	Now, we divide the total feed needed
1093		(3):	for the day by the number of meals
1094		800 cups / 3 meals = 266.67 cups	(3) and the number of chickens (20):
1095		So, Wendi needs to give her chickens	800  cups / (3  meals  * 20  chickens) =
1096		200.0 / cups of feed in the final meal of the day	$\delta 00 \text{ cups } / (\delta 0 \text{ meals } \circ \text{ chickens}) =$
1097		or the day.	So. Wendi needs to give her chickens
1098			13.33 cups of feed in the final meal
1099			of the day.
1100	Question: Kylar went to the store to	Kylar needs to pay \$80 for 16 glasses.	Kylar needs to pay \$80 for 16 glasses.
1101	buy glasses for his new apartment.	Confidence: 95%	Confidence: 100%
1102	One glass costs $$5$ , but every second		
1103	lar wants to buy 16 glasses How		
1104	much does he need to pay for them?		
1105	Question: Toulouse has twice	To find the total number of sheep	Toulouse has twice as many sheep
1106	as many sheep as Charleston.	in Toulouse, Charleston, and Seat-	as Charleston, so if Charleston has
1107	Charleston has 4 times as many	tle together, we need to first find the	4 times as many sheep as Seattle,
1108	sheep as Seattle. How many sheep	number of sheep in Charleston and	Toulouse must have 8 times as many
1109	do Ioulouse, Charleston, and Seattle	Toulouse. Since Charleston has 4	sheep as Seattle. Since Seattle has 20 sheep. Toulouse has 160 sheep.
1110	have together if Seattle has 20 sheep.	Seattle has 20 sheep. Charleston has	Now, let's find the total number of
1111		$4 \ge 20 = 80$ sheep. Now, we know	sheep for Toulouse, Charleston, and
1112		that Toulouse has twice as many	Seattle together:
1113		sheep as Charleston, so Toulouse has	Toulouse: 160 sheep Charleston:
1114		$2 \ge 80 = 160$ sheep.	4 times as many sheep as Seattle, which is $4 \times 20 = 80$ sheep Seattle;
1115		in Toulouse. Charleston, and Seattle	20 = 80  sheep Seattle.
1116		together:	Total sheep: $160 + 80 + 20 = 260$
1117		Toulouse + Charleston + Seattle =	sheep
1118		160 + 80 + 20 = 260 sheep	So, the total number of sheep for
1119		So, the total number of sheep in Toulouse Charleston and Soattle to	Toulouse, Charleston, and Seattle to-
1120		gether is 260.	geulei 15 200.
1121		Search 15 200.	
1122			
1123			
1124			
1125			
1126			
1127			
1128			
1120			
1130			
1121			
1131			

1134			
1135	Question: Carla is downloading a	Carla's download speed is 2 CP/minute The file size is 200 CP	Carla's download speed is 2 GP/minute The file size is 200 GP
1136	load 2 GB/minute but 40% of the	The download will take 100 minutes	The download will take 100 minutes
1137	way through the download, Windows	(200 GB / 2 GB/minute) to complete.	(200 GB / 2 GB/minute) to complete.
1138	forces a restart to install updates,	However, 40% of the way through the download, Windows forces a restart to install updates, which takes 20 minutes. After the restart, Carla has to restart the download from	However, 40% of the way through
1139	which takes 20 minutes. Then Carla		the download, Windows forces a
1140	has to restart the download from the		restart to install updates, which takes
1141	download the file?		has to restart the download from the
1142		the beginning. This means that the	beginning.
1143		download will take an additional 80	So, the total time to download the
1144		minutes (200 GB / 2 GB/minute) to	file is 120 minutes (100 minutes +
1145		required to download the file is 100	20 minutes). The download speed is
1146		minutes $+ 20$ minutes $+ 80$ minutes	utes).
1147		= 200 minutes.	
1148	Question: John drives for 3 hours at	John is 120 miles from home at the	John's total distance traveled can be
1149	a speed of 60 mph and then turns	end of the 4 hours.	calculated by adding the distances
1150	around because he realizes he forgot	Confidence: 95%	covered during each leg of his jour-
1151	He tries to get home in 4 hours but		mph which gives a distance of 3 x
1152	spends the first 2 hours in standstill		60 = 180 miles. The second leg is 4
1153	traffic. He spends the next half-hour		hours, with the first 2 hours spent in
1154	driving at a speed of 30mph, before		standstill traffic, the next half-hour at
1155	being able to drive the remaining		30 mph, and the remaining time at 80
1156	How far is he from home at the end		is:
1157	of those 4 hours?		2 hours at 0 mph (standstill traffic) = $\frac{1}{2}$
1158			0 miles 0.5 hours at 30 mph = $0.5 \text{ x}$
1159			30 = 15 miles 1.5 hours at 80 mph =
1160			$1.5 \times 80 = 120 \text{ miles}$
1161			135 = 315 miles. Since John turned
1162			around after driving 180 miles, he is
1163			135 miles from home at the end of
1164			the 4 hours.
1165	Question: Eliza's rate per hour for the	Eliza's earnings for this week can be	Eliza's earnings for this week can be
1166	first 40 hours she works each week	calculated as follows:	calculated as follows:
1167	pay of 1.2 times her regular hourly	first 40 hours: - Regular hourly rate:	first 40 hours: - Regular hourly rate:
1168	rate. If Eliza worked for 45 hours this	\$10 - Hours worked: 40 - Regular	\$10 - Hours worked: 40 - Regular
1169	week, how much are her earnings for	pay: \$10 x 40 = \$400	pay: \$10 x 40 = \$400
1170	this week?	2. Calculate the overtime pay for the	2. Calculate the overtime pay for the
1171		additional 5 nours: - Overtime hourly rate: $1.2 \times $10 - $12$ - Hours worked:	remaining 5 nours: - Overtime hourly rate: $1.2 \times $10 - $12$ - Hours worked:
1172		5 - Overtime pay: $$12 \times 5 = $60$	5 - Overtime pay: $$12 \times 5 = $60$
1173		3. Add the regular pay and overtime	3. Add the regular pay and overtime
1174		pay to get the total earnings: - Total	pay to get the total earnings: - Total
1175		earnings: \$400 + \$60 = \$460	earnings: $$400 + $60 = $460$
1176		\$460	SO, Eliza's earnings for this week are
1177		•	
1178			
1179			

1188 1189 1190 1191 1192 1193 1194 1195 1196 1197 1198 1199 1200 1201 1202	Question: A new program had 60 downloads in the first month. The number of downloads in the second month was three times as many as the downloads in the first month, but then reduced by 30% in the third month. How many downloads did the program have total over the three months?	In the first month, the program had 60 downloads. In the second month, the number of downloads was three times as many as the downloads in the first month, so the number of downloads in the second month is $3 * 60 = 180$ . In the third month, the number of downloads reduced by 30%, so the number of downloads in the third month is $180 * 0.7 = 126$ . The total number of downloads over the three months is $60 + 180 + 126 = 366$ .	In the first month, the program had 60 downloads. In the second month, the number of downloads was three times as many as the downloads in the first month, so the number of downloads in the second month is $3 * 60 = 180$ . In the third month, the number of downloads reduced by 30% from the downloads in the second month, so the number of downloads in the second month, so the number of downloads in the third month is $0.7 * 180 = 126$ . The total number of downloads over the three months is: 60 + 180 + 126 = 366 So, the program had a total of 366 downloads over the three months.
1202 1203 1204 1205 1206 1207 1208 1209 1210 1211 1212 1213 1214 1215	Question: Toula went to the bakery and bought various types of pastries. She bought 3 dozen donuts which cost \$68 per dozen, 2 dozen mini cup- cakes which cost \$80 per dozen, and 6 dozen mini cheesecakes for \$55 per dozen. How much was the total cost?	To calculate the total cost, we need to add the costs of the different pastries. The cost of 3 dozen donuts is \$68 per dozen, so the total cost is \$68 * 3 = \$204. The cost of 2 dozen mini cup- cakes is \$80 per dozen, so the total cost is \$80 * 2 = \$160. The cost of 6 dozen mini cheesecakes is \$55 per dozen, so the total cost is \$55 * 6 = \$330. Now, we can add all the costs together to get the total cost: \$204 + \$160 + \$330 = \$694 So, the total cost of the pastries Toula bought was \$694.	Toula's total cost for the pastries can be calculated by adding the costs of the donuts, mini cupcakes, and mini cheesecakes. The cost of the donuts is \$68 per dozen, so the total cost for 3 dozen is \$68 * 3 = \$204. The cost of the mini cupcakes is \$80 per dozen, so the total cost for 2 dozen is 80 * 2 = \$160. The cost of the mini cheesecakes is \$55 per dozen, so the total cost for 6 dozen is \$55 * 6 = \$330. Adding these costs together, we get: \$204 + \$160 + \$330 = \$694 So, the total cost of Toula's pastries was \$694.
1216 1217 1218 1219 1220 1221 1222 1223 1224 1225 1226 1227 1228 1229 1230 1231 1232 1233 1234 1235 1236 1237 1238	Question: Carlos is planting a lemon tree. The tree will cost \$90 to plant. Each year it will grow 7 lemons, which he can sell for \$1.5 each. It costs \$3 a year to water and feed the tree. How many years will it take be- fore he starts earning money on the lemon tree?	To determine the number of years it will take for Carlos to start earning money on the lemon tree, we need to calculate the break-even point. The break-even point is the point at which the total revenue from selling lemons equals the total cost of planting and maintaining the tree. The total cost of planting the tree is \$90. The annual cost of watering and feeding the tree is \$3. The annual rev- enue from selling lemons is 7 lemons multiplied by \$1.5 per lemon, which is \$10.5. We can calculate the break- even point by setting the total cost equal to the total revenue: \$90 + \$3 * (years) = \$10.5 * (years) Solving for years, we get: years = 90 + 0.2857 * (years) years = 90 + 0.2857 * (years - 90) years = 90 + 0.2857 * years - 25.71 0.7143 * years = 64.29 years = 90.1 So, it will take approximately	Carlos will start earning money on the lemon tree after the first year. Here's why: 1. The tree costs \$90 to plant. 2. Each year, the tree will grow 7 lemons, which he can sell for \$1.5 each. 3. It costs \$3 a year to water and feed the tree. After the first year, Carlos will have earned \$10.5 from selling the lemons (\$1.5 x 7 lemons) and spent \$3 on maintenance. This means he will have made a profit of \$7.5 (\$10.5 - \$3) after the first year. So, he will start earning money on the lemon tree after the first year.

1242	Question: Melanie is a door to door	Melonia started with 15 vacuum	Malania started with 15 vacuum
1243	saleswoman. She sold a third of her	cleaners.	cleaners.
1244	vacuum cleaners at the green house,		
1245	2 more to the red house, and half of		
1246	what was left at the orange house. If		
1247	how many did she start with?		
1248	Ouestion: In a dance class of 20 stu-	In a dance class of 20 students, 20%	In a dance class of 20 students, 20%
1249	dents, 20% enrolled in contemporary	enrolled in contemporary dance, 25%	enrolled in contemporary dance, 25%
1250	dance, 25% of the remaining enrolled	of the remaining enrolled in jazz	of the remaining enrolled in jazz
1251	in jazz dance, and the rest enrolled	dance, and the rest enrolled in hip-	dance, and the rest enrolled in hip-
1252	of the entire students enrolled in hip-	To find the percentage of students en-	To find the percentage of students en-
1253	hop dance?	rolled in hip-hop dance, we need to	rolled in hip-hop dance, we need to
1254	-	calculate the number of students en-	calculate the number of students en-
1255		rolled in hip-hop dance and divide it	rolled in hip-hop dance and divide it
1256		by the total number of students in the	by the total number of students in the
1257		First, let's find the number of students	First, let's find the number of students
1200		enrolled in contemporary dance:	enrolled in contemporary dance:
1260		20% of 20 students = $20% * 20 = 0.2$	20% of 20 students = $20% * 20 = 0.2$
1261		Now, let's find the number of stu- dents enrolled in jazz dance:	Now, let's find the number of stu- dents enrolled in jazz dance:
1262			
1263		25% of the remaining students = $25\%$	25% of the remaining students = $25\%$
1264		*(20 - 4) = 25% * 16 = 0.25 * 16 = 4	*(20 - 4) = 25% * 16 = 0.25 * 16 = 4
1265		students The number of students enrolled in	students The number of students enrolled in
1266		hip-hop dance is the remaining 12	hip-hop dance is the remaining 12
1267		students:	students:
1268		Total students - (contemporary dance	Total students - (contemporary dance
1269		+ jazz dance) = $20 - (4 + 4) = 20 - 8$ - 12 students	+ jazz dance) = $20 - (4 + 4) = 12$
1270		Now, let's find the percentage of stu-	Now, let's find the percentage of stu-
1271		dents enrolled in hip-hop dance:	dents enrolled in hip-hop dance:
1272		Percentage of students enrolled in	Percentage of students enrolled in
1273		hip-hop dance = $(12 / 20) * 100 =$	hip-hop dance = $(12 / 20) * 100 =$
1274		So. 60%	00 %
1275	Ouestion: A merchant wants to make	To maximize profit, the merchant	To maximize profit, the merchant
1276	a choice of purchase between 2 pur-	should choose the electronic gadgets	should choose the electronic gadgets
1277	chase plans: jewelry worth \$5,000	worth \$8,000. Assuming the finan-	worth \$8,000. Assuming the finan-
1278	or electronic gadgets worth \$8,000.	cial advisor's speculation is correct,	cial advisor's prediction is accurate,
1279	the jewelry market will go up 2.5%	and the electronic gadgets market	Profit = (1.12 * \$8.000) - \$8.000 =
1280	while the electronic gadgets market	will rise 1.2% within the same month.	\$896
1281	will rise 1.2% within the same month.	The profit from the jewelry purchase	The merchant would make a profit
1282	If the merchant is looking to maxi-	Would be: Profit from jewelry $-(5.000 * 2.5\%)$	of \$896 by choosing the electronic
1283	by making a choice, how much profit	= 125	gauguis.
1284	would this be?	The profit from the electronic gadgets	
1200		purchase would be:	
1200		Profit from electronic gadgets = $(8,000 \times 1,2\%) = 06$	
120/		$(0,000 \cdot 1.2\%) = 90$ Since the profit from the electronic	
1200		gadgets purchase is higher, the mer-	
1203		chant should choose the electronic	
1291		gadgets worth \$8,000 to maximize	
1292		piont.	
1293			

1296 1297 1298 1299 1300 1301	Question: Two trains leave San Rafael at the same time. They be- gin traveling westward, both trav- eling for 80 miles. The next day, they travel northwards, covering 150 miles. What's the distance covered by each train in the two days?	In the first day, both trains travel 80 miles westward. In the second day, they travel 150 miles northward. The total distance covered by each train in the two days is 230 miles.	In the first day, both trains travel 80 miles westward. In the second day, they travel 150 miles northward. The total distance covered by each train in the two days is 230 miles.
1302 1303 1304 1305 1306 1307 1308 1309 1310 1311 1312 1313 1314 1315 1316 1317	Question: Jill gets paid \$20 per hour to teach and \$30 to be a cheerleading coach. If she works 50 weeks a year, 35 hours a week as a teacher and 15 hours a week as a coach, what's her annual salary?	Jill's annual salary can be calculated by multiplying her hourly rate, the number of hours she works per week, and the number of weeks she works in a year. Jill's hourly rate for teaching is \$20, and she works 35 hours per week. Her weekly salary for teaching is: \$20 imes 35 = \$700 Jill's hourly rate for coaching is \$30, and she works 15 hours per week. Her weekly salary for coaching is: \$30 imes 15 = \$450 Jill works 50 weeks a year, so her annual salary is: \$700 imes 50 + \$450 imes 50 = \$35,000 + \$22,500 = \$57,500 So, Jill's annual salary is \$57,500.	Jill's annual salary can be calculated by multiplying her hourly rate, the number of hours she works per week, and the number of weeks she works in a year. Here's the calculation: Teacher's annual salary: \$20 per hour x 35 hours per week x 50 weeks = \$35,000 Coach's annual salary: \$30 per hour x 15 hours per week x 50 weeks = \$22,500 Total annual salary: \$35,000 + \$22,500 = \$57,500 So, Jill's annual salary is \$57,500.
1317 1318 1319 1320 1321 1322 1323 1324 1325 1326 1327 1328	Question: Claire makes a 3 egg omelet every morning for breakfast. How many dozens of eggs will she eat in 4 weeks?	In 4 weeks, Claire will eat 12 dozens of eggs if she makes a 3 egg omelet every morning. This is because there are 4 weeks in a month, and each month has 4 weeks. So, in 4 weeks, she will make 12 omelets, and each omelet requires 3 eggs. Therefore, she will consume 12 dozens of eggs in 4 weeks.	In 4 weeks, Claire will eat 48 dozens of eggs if she makes a 3 egg omelet every morning. This is because there are 4 weeks in a month, and each week has 7 days. If she makes a 3 egg omelet every day, she will consume 21 eggs per week. 21 eggs multiplied by 4 weeks equals 84 eggs. Since there are 12 eggs in a dozen, 84 eggs divided by 12 equals 7 dozens. There- fore, Claire will eat 7 dozens of eggs per week, and 7 dozens multiplied by 4 weeks equals 28 dozens.
1329 1330 1331 1332 1333 1333	Question: Marissa is hiking a 12- mile trail. She took 1 hour to walk the first 4 miles, then another hour to walk the next two miles. If she wants her average speed to be 4 miles per hour, what speed (in miles per hour) does she need to walk the remaining distance?	To maintain an average speed of 4 miles per hour, Marissa needs to walk the remaining distance at a speed of 4 miles per hour.	To maintain an average speed of 4 miles per hour, Marissa needs to walk the remaining distance at a speed of 4 miles per hour.
1335 1336 1337 1338 1339 1340 1341 1342 1343 1344 1345 1346 1347 1348 1349		1	1