# EVERYTHING EVERYWHERE ALL AT ONCE: LLMS CAN IN-CONTEXT LEARN MULTIPLE TASKS IN SUPERPOSI TION

Anonymous authors

Paper under double-blind review

### ABSTRACT

Large Language Models (LLMs) have demonstrated remarkable in-context learning (ICL) capabilities. In this study, we explore a surprising phenomenon related to ICL: LLMs can perform multiple, computationally distinct ICL tasks simultaneously, during a single inference call, a capability we term "task superposition". We provide empirical evidence of this phenomenon across various LLM families and scales and show that this phenomenon emerges even if we train the model to in-context learn one task at a time. We offer theoretical explanations that this capability is well within the expressive power of transformers. We also explore how LLMs internally compose task vectors during superposition. Furthermore, we show that larger models can solve more ICL tasks in parallel, and better calibrate their output distribution. Our findings offer insights into the latent capabilities of LLMs, further substantiate the perspective of "LLMs as superposition of simulators", and raise questions about the mechanisms enabling simultaneous task execution.

025 026

005 006

007

008 009 010

011

013

014

015

016

017

018

019

021

023

## 1 INTRODUCTION

027 028

052

Large Language Models (LLMs) have demonstrated remarkable capabilities across various domains,
with one of the most intriguing being in-context learning (ICL) (Brown et al., 2020; Xie et al., 2022).
ICL enables LLMs to perform tasks during inference without the need to fine-tune for that particular
task, simply by providing a few examples within the input prompt. This ability has sparked significant
interest in the research community, as it suggests that LLMs can adapt to novel tasks on-the-fly, using
the capabilities that they acquired during pretraining, and the context provided.

While ICL has been extensively studied from both theoretical and empirical perspectives (Xie et al., 2022; Agarwal et al., 2024), many aspects of its underlying mechanisms remain elusive. In this work, we study a surprising phenomenon related to ICL that, to the best of our knowledge, has not been thoroughly studied before: LLMs can perform multiple distinct ICL tasks simultaneously, in a single inference call, a capability we refer to as *"task superposition"*.

Our study suggests that pretrained autoregressive LLMs such as Llama (Touvron et al., 2023) or GPT-3.5<sup>1</sup> (Brown et al., 2020) display superposition of tasks purely *in-context*<sup>2</sup>. When presented with multiple in-context examples from different tasks, in the same prompt, the models can generate outputs that correspond to solutions for all these individual tasks. For instance, given examples of addition and translation, the model can concurrently produce correct answers for both tasks, as well as the composition of these tasks (e.g., the result of addition translated into another language).

Figure 1 illustrates this phenomenon. In Figure 1a (left), given in-context examples of addition in different languages and the query "91 + 83  $\rightarrow$ ", the model generates probabilities for the correct sum in various languages, demonstrating its ability to perform addition and translation concurrently.

This discovery aligns and lends further support to the view of LLMs as superposition of simulators (Janus, 2022; Shanahan et al., 2023; Nardo, 2023) and the Bayesian perspective of ICL proposed by Xie et al. (2022). While not a mathematically rigorous formulation, we can conceptualize the output

<sup>&</sup>lt;sup>1</sup>In particular, gpt-3.5-turbo-instruct.

<sup>&</sup>lt;sup>2</sup>For other definitions of superposition, please see Section 2.

069

071

073

075

076

077

078

079

081

082

083 084

085

087

090

091 092

094



(a) (left) Two-digit addition in a variety of languages. (right) Naming the capital of a given country name, naming the continent of a given country name or capitalizing the country name.



(b) (left) Tasks copy (op1), copy (op2) and op1+op2. (right) First or last letter in upper or lower case.

Figure 1: LLMs can perform task superposition. (a) Llama-3 70B and (b) GPT-3.5 Turbo are each presented with two sets of tasks. For each set of tasks, we show an example prompt such that all except the last row are in-context examples of one of the tasks and the last row is the query. We provide 20 in-context task examples for each task in the prompt with order randomized and provide the probabilities of outputs when correctly performing each task on the query.

of an LLM as a weighted sum of conditional probabilities across possible tasks:

$$\mathbb{P}(\mathsf{output}|\mathsf{prompt}) \approx \sum_{\mathsf{task}} \mathbb{P}(\mathsf{output}|\mathsf{task},\mathsf{prompt}) \mathbb{P}(\mathsf{task}|\mathsf{prompt}).$$

In this conceptual model,  $\mathbb{P}(\text{output}|\text{prompt})$  represents the probability distribution over possi-095 ble outputs given the input prompt, a task can be thought of as a latent variable representing 096 different capabilities the model might possess (e.g., arithmetic, translation, sentiment analysis),  $\mathbb{P}(\text{output}|\text{task}, \text{prompt})$  represents the output probability distribution if the model was specifically 098 attempting to solve a single task, based on the test example in the prompt, and  $\mathbb{P}(\mathsf{task}|\mathsf{prompt})$ 099 represents the model's inferred probability that the prompt specifies a particular task. 100

While this mental model is a simplification of how an LLM operates, it provides an intuitive way to 101 support the task superposition phenomenon we observe. Our findings lend support to the idea that 102 LLMs can simultaneously maintain and utilize multiple task distributions, resulting in outputs that 103 reflect a combination of relevant tasks. 104

# 105 106

## **Our Contributions:** Our study makes several key contributions:

1. Through extensive empirical investigation and theoretical results, we demonstrate that task 107 superposition is prevalent across various pretrained LLM families (GPT3.5, LLama-3, Qwen).

- 108 2. We empirically show that task superposition emerges as we train on one task at a time.
  - 3. We provide a theoretical construction showing that Transformers models are indeed capable of task superposition, and have the capacity to implement multiple tasks in parallel.
  - 4. We explore how LLMs internally compose task vectors (Hendel et al., 2023) during superposition, and show how convex combinations of task vectors can reproduce the superposition effect.
  - 5. We show that larger models can solve more tasks in parallel and more accurately reflect the distribution of in-context tasks.

We believe that our findings offer new insights into the latent capabilities of LLMs and raise questions about the mechanisms enabling simultaneous task execution. We believe this work sheds more light on the ICL capabilities of frontier language models, and offers a glimpse on potential applications of task superposition in practical settings.

119 120 121

110

111

112

113

114

115 116

117

118

# 2 RELATED WORK

122 123

Theory and practice of in-context learning. There is rich literature which formalizes in-context 124 learning under diverse definitions. For example, prior works study in-context learning through a 125 Bayesian framework for task retrieval (Xie et al., 2022; Panwar et al., 2023; Zhang et al., 2023), 126 martingales (Falck et al., 2024), optimizers (Akyürek et al., 2023; Oswald et al., 2023; Dai et al., 127 2022) and more (Reddy, 2024; Olsson et al., 2022). Other works confirm the theoretical framing of 128 in-context learning by using it to implement a variety of algorithms and methods (Zhou et al., 2023; 129 Ahn et al., 2023; Giannou et al., 2023; Wu et al., 2024; Laskin et al., 2022; Zhou et al., 2022), or to 130 approximate general-purpose computing machines (Giannou et al., 2023; Wei et al., 2022). 131

To bridge the gap between theory and practice, many works have used these theoretical insights to study in-context learning behaviors, such as in many-shot in-context learning, (Agarwal et al., 2024), long-context (Li et al., 2024), or eliciting personas (Choi & Li, 2024). Other works study the factors that influence how well models can learn through context, such as task diversity (Raventos et al., 2023; Chan et al., 2022), the balance between pre-training priors and in-context (Wei et al., 2023; Lin & Lee, 2024), in-context labels (Min et al., 2022; Lyu et al., 2022), and the in-context format (Lampinen et al., 2022). In-context learning has also been proposed as a means of fine-tuning to improve non-language tasks (Dinh et al., 2022).

The development of new architectures such as state space models (Gu & Dao, 2023) has further motivated studying whether in-context learning is prevalent in alternative architectures such as Mamba (Park et al., 2024; Grazzi et al., 2024; Zeng et al., 2024) or in looped transformers (Yang et al., 2023).

Steering models through in-context learning has been a growing area of interest. Recent work has hypothesized that in-context learning can be encapsulated by a high-dimensional description of a task, which can be used to replace, (Hendel et al., 2023) compose (Todd et al., 2024) or augment (Liu et al., 2024) the latent states of a model, in order to alter its default behavior. Task vectors can be combined via arithmetic operations to solve a variety of tasks (Ilharco et al., 2023). Prior work has also been investigating the power of tokens in defining a task (Bai et al., 2024).

Other definitions of superposition. Our findings on superposition are inspired by notions of language
 models as multiverse generators (Reynolds & McDonell, 2021; moire, 2021). One consequence of
 LLMs as a superposition of tasks is that the outputs may collapse to unintended simulacra, a behavior
 known as the "Waluigi effect" (Nardo, 2023).

Superposition has been defined in various related contexts of learning models. Feature superposition
(Elhage et al., 2022) refers to a neural network's ability to represent multiple learned concept in
a single neuron. Though our discovery of task superposition describes the same abstract idea, we
stress that it is distinct from feature superposition because task superposition is most apparent in
the final output of a model. Feature superposition is a microscopic-level observation whereas task
superposition is a macroscopic-level observation.

- Superposition is also described as a way to store multiple models in a single set of parameters (Cheung et al., 2019), processing multiple inputs simultaneously (Shen et al., 2024a; Murahari et al., 2022). In our work, we demonstrate task superposition directly as a result of language pre-training,
- without the necessity of additional adapters or decoding strategies.



Figure 2: Distributions (0/25/50/75/100-percentiles) of probabilities for correct outputs of each task. For every set of tasks, we tested with 100 prompts and for each prompt, every task has 20 random in-context task examples with order randomized like in Figure 1. Category other is the sum of probabilities of all other outputs. Gray dashed line in each figure is the ideal probability if we assume the model perfectly calibrates its output distribution to the distribution of in-context task examples. With uniform distribution of task examples, the dashed lines are at 0.25 (4 tasks setting) and 0.33 (3 tasks setting).

206 207

208

209

210 211 212

# 3 LLMS ARE A SUPERPOSITION OF MULTIPLE IN-CONTEXT LEARNERS

In this section, we want to investigate if existing pre-trained models exhibit superposition of multiple tasks and whether this phenomenon is common (i.e., whether we can observe this phenomenon on a variety set of tasks and different families of LLMs).

**Finding 1:** *LLMs can in-context learn multiple tasks in superposition when provided with prompts of a mixture of task examples.* 

213 214 215

We denote K by the number of tasks and consider four different settings of task mixtures.

1. Numerical addition and addition in English, French or Spanish (K = 4). Example prompt is shown in Figure 1a (left).

- 218 2. Given a name of a country, name the capital, continent or capitalize the country name  $(V_{i}, 2)$ . Ensure the ground is the capital continent of capitalize the country name
- 220

216

217

221

(K = 3). Example prompt is shown in 1a (right).
3. Given input "{op1}@{op2}", copy op1, op2 or add op1 and op2 (K = 3). Example prompt is shown in Figure 1b (left).

222 223 224

4. Given a word, output first letter or last letter in lower or upper cases (K = 4). Example prompt is shown in 1b (right).

We provide GPT-3.5 (Brown et al., 2020), Llama-3 70B (AI@Meta, 2024) and Qwen-1.5 72B (Bai et al., 2023a) with prompts of uniform mixture of tasks (each task has 20 examples in the prompt ordered randomly). For each prompt consisting of in-context task examples (e.g., " $11 + 26 \rightarrow 37$ " for the first task in the first setting) and a query (e.g., " $91 + 83 \rightarrow$ "), we calculate the probabilities of outputs when correctly performing each task on the query and plot the distribution of probabilities for each task in Figure 2. Details on calculating the probabilities is in Appendix B.

Figure 2 reveals that in all four sets of tasks, all models have non-negligible median values of probabilities for at least two tasks. This indicates that the models can in-context learn multiple tasks in superposition when provided with prompts of a mixture of task examples.

We can also observe that, even though every task in a prompt has an equal number of in-context examples (20 examples), LLMs do not calibrate their output distribution perfectly with the in-context task example distribution and they still have bias on what task to perform. For example, Figure 2a shows that Llama-3 70B prefers performing numerical addition over addition in other languages, Qwen-1.5 72B prefers addition in English while GPT-3.5 does not have a strong preference over a single task. On the other hand, in Figure 2b GPT-3.5 has a strong preference over the capital task.

Additionally, some tasks are "harder" than other tasks. For example, in Figure 2d, all models assign near-zero probability for task answers of last\_letter and last\_letter\_cap. The category other has relatively high values, indicating a high noise when prompted with in-context examples of this setting. In contrast, in Figure 2c, category other has very small values, indicating that all models most of the time would correctly assign the output probabilities to the correct answers.

245 246 247

# 4 TASK SUPERPOSITION IN MODELS TRAINED FROM SCRATCH

In Section 3 we investigated task superposition in pre-trained LLMs at inference time. We further
 investigate how task superposition emerges in LLMs during training. Specifically, if we train the
 model to in-context learn one task at a time, can it perform task superposition when provided with
 prompts containing examples of multiple tasks?

252 To answer this question, we train a small GPT-2 model (6 heads, 6 layers,  $\sim$ 14million pa-253 rameters) (Radford et al., 2019) to learn a family of retrieval tasks. The input has the form 254 "{ch1}{ch2}{ch3}{ch4}{ch5}{ch6}{ch7}{ch8}→" where ch1, ..., ch8 are distinct sin-255 gle characters. We consider 8 retrieval tasks - ret1, ..., ret8 - where ret1 is to output ch1 256 and so on. The model is trained to in-context learn one task (retrieve one of {ch1, ..., ch8}) at a 257 time in training. Namely, during training, the model is only provided with text data such that each 258 prompt only contains in-context examples of a single randomly chosen task (and different prompts 259 can correspond to different tasks).

Concretely, for each sample, we randomly select task  $t \in \{\text{ret1}, ..., \text{ret8}\}$  and inputs  $x^{(1)}, ..., x^{(m)}$ , where each  $x^{(j)}$  is an eight-character long string. We then form the sequence  $s = [x^{(1)}, g_t(x^{(1)}), ..., x^{(m)}, g_t(x^{(m)})]$  where  $g_t(x^{(j)})$  is the output of performing task t on  $x^{(j)}$ . We train the model  $M_\theta$  parametrized by  $\theta$  using ICL training. In particular, we minimize the following objective:

265 266 267

$$\min_{\theta} \mathbf{E}_{\boldsymbol{s}} \left( \frac{1}{m-1} \sum_{j=1}^{m-1} \operatorname{CE}(M_{\theta}(\boldsymbol{s}_{j} \oplus \boldsymbol{x}^{(j+1)}), \boldsymbol{g}_{t}(\boldsymbol{x}^{(j+1)})) \right),$$
(1)

where 
$$s_j \oplus x^{(j+1)} \equiv [x^{(1)}, g_t(x^{(1)}), ..., x^{(j)}, g_t(x^{(j)}), x^{(j+1)}]$$
 and CE is the cross-entropy loss.







Figure 3: We consider two different settings of tasks: (a) given an eight-character length string as input, consider ret1, ..., ret8 where ret1 is to retrieve the first character and so on; and (b) 288 given a two-digit integer as an input, consider plus0, ..., plus9 where plus0 is to add 0 on the 289 input and so on. After training, for each setting, we select two tasks ad we provide the model with 290 prompts containing in-context examples from these two tasks and vary the mixture ratio  $\lambda$  such that the in-context task example distribution for two tasks is  $[\lambda, 1-\lambda]$ . We plot  $\lambda$  on x-axis and the output 292 probabilities of task answers for each task on y-axis. 293

After training, we provide the model with prompts containing in-context examples of two tasks (in particular, we choose ret2 and ret6) and see if the model performs task superposition. We vary the proportion of in-context examples of two tasks and plot the output distributions in Figure 3a.

Similarly, we consider a second setting involving 10 tasks. Given a two digit integer input num, task plus0 outputs num, task plus1 outputs num + 1 and so on, up to task plus9. The model is trained to in-context learn one of plus0,..., plus9 at a time, following the procedure above. During inference time, the model is tested with prompts containing a mixture of in-context examples from tasks plus2 and plus6. We vary the mixture ratio and show the output distributions in Figure 3b.

Finding 2: Transformers can in-context learn multiple tasks in superposition even if trained to in-context learn one task at a time.

Remarkably, from Figure 3a and 3b, GPT-2 trained from scratch to in-context learn one task at a time can generalize to simultaneously performing multiple tasks and calibrate the output probabilities according to the in-context task example distribution when provided with a mixture of in-context examples. For example, in Figure 3a at the mixture ratio  $\lambda = 0.5$ , meaning that 50 percent of the examples in the prompt is from task ret2 and the other 50% comes from task ret6, we can see the output probabilities for task answers of ret2 and ret6 being roughly [0.5, 0.5]. We can observe similar behavior in Figure 3b.

313 314

284

285

286

287

291

295

296

297

298

299

300

301

302 303 304

305 306 307

308

309

310

311

312

315 316

5

## TRANSFORMERS HAVE THE CAPACITY TO PERFORM TASK SUPERPOSITION

317 In this section, we explore whether Transformers have the inherent expressivity to perform multiple 318 tasks in superposition with a single inference call. To this end, we provide a theoretical construction 319 of a Transformer which, given the ability to implement multiple tasks, performs task superposition 320 depending on the examples given in-context. 321

**Theorem 1.** A seven layer transformer with embedding dimension  $\mathcal{O}(d + \log(mn))$  with K heads 322 per attention layer can perform K tasks on vectors of dimension d in superposition, with weighting 323 based on m different in-context examples each of length n.

The proof of Theorem 1 is provided in Appendix D.4. Note that while this does not guarantee that training a Transformer will actually find these parameters, it does indicate that Transformers are expressive enough to perform task superposition at test time. Below we outline the main ideas used in the proof.

**Prediction based on multiple tasks.** Assume that we are given m in-context samples  $(x_1^{(j)}, \ldots, x_{n-2}^{(j)}, \stackrel{=}{=}, y^{(j)})_{j=1}^m$  where '=' represents a specific value used only for preceding the label, and a set of k different Transformers TF<sub>i</sub> which can implement the T different desired tasks, where each deterministic task is denoted as  $g_i(x^{(j)})$  with  $i \in [k]$  and  $j \in [m]$ , *i.e.*  $y^{(j)} = g_i(x^{(j)})$ for some task i dependent on sample j. Using the weights of each TF<sub>i</sub>, we can compute outputs of the following form:

[	$x_1^{(j)}$		$x_{n-2}^{(j)}$	=	$oldsymbol{y}^{(j)}$	]	[···	$x_1^{(j)}$		$oldsymbol{x}_{n-2}^{(j)}$	=	$\boldsymbol{y}^{(j)}$	]
	0	• • •	0	0	0	···  _	$\rightarrow$  ···	0	•••	0	0	$\ m{g}_1(m{x}^{(j)}) - m{y}^{(j)}\ _1$	
İ	:		:	:	:	i		:		:	:	:	İ
[	0		0	0	0	]	[	0		0	0	$\ m{g}_T(m{x}^{(j)}) - m{y}^{(j)}\ _1$	]

We use the  $l_1$  norm to aggregate the prediction, in case that the task is multi-dimensional. These differences are used to identify tasks, as  $\|g_i(x^{(j)}) - y^{(j)}\|_1 \approx 0$  for  $y^{(j)}$  coming from task *i*. Different heads at each layer in the model are used to execute each of the tasks in parallel using the weights from TF<sub>i</sub>. In Appendix D we construct tasks where an arbitrary function  $g_i(x_l^{(j)})$  is implemented using ReLUs for some fixed *l* that is task-specific.

346 **Creating task identifiers.** Having the differences between the implemented function and the 347 label, we first use the ReLUs to clean up the vectors  $v_k$  so that only the positions in each vector 348 that are associated with a task are maintained and the rest are set to 1<sup>3</sup>. We thus create the vectors  $(v'_k)_i = \|g_k(x_{1:l-1}^{(j)}) - x_l^{(j)}\|_1$  and  $(v'_*)_* = 1$  otherwise. Now we use ReLUs to threshold and create an indicator vectors  $\mathbf{1}_{\{\|g_k(x_{1:l-1}^{(j)}) - x_l^{(j)}\|_1 \approx 0\}}$  which identify the task, *i.e.*,these are task identifiers. 349 350 351 Notice that if the task is correctly predicted then the difference should be close to 0 (up to some 352 error), while if the task is not identified the corresponding value would not be 0; the rest of the 353 rows would be 1. We have created one vector for each task, which has 1 in the position of the 354 corresponding task if the task was identified in the context. 355

Averaging and task superposition. As a last step, we average all the task identifiers and place the
 result in the last column, in which the next prediction will happen. We then use the averaged task
 identifier to weight the prediction of each task based on it, as in task superposition. If the task has
 been identified multiple times in the context, it would be assigned a higher weight/probability.

360 361 362

364

365

366

367

368

369

373

374

375

376

377

# 6 TASK SUPERPOSITION THROUGH THE LENS OF TASK VECTORS

While in Section 5 we provide an existential result by constructing a Transformer that performs task superposition and shows that task superposition is well within the expressive power of Transformers, we would like to further investigate how task superposition manifest in pretrained LLMs internally. In this section we explore the underlying mechanisms that LLMs employ during task superposition. In particular, we focus our empirical study on *task vectors* (Hendel et al., 2023) where the detailed implementation is in Appendix C. Task vectors are vectors in the embedding space and are found to encode the algorithm that a model internally implements to solve a task given in-context demonstrations.

We want to investigate if there is any relation between the task vectors of each individual task and the task vectors of a mixture of task examples in the prompt. To this end, we consider two sets of tasks:

- (a) copy (op1), copy (op2) and op1+op2 as in Figure 1b (left).
- (b) Given a two-digit integer, task to\_fr translates it to French, task to\_de translates it to German and task to\_it translates it to Italian.

<sup>&</sup>lt;sup>3</sup>This step is not mandatory, but it ensures that we have no values over which we have no control. We leave as future work an error analysis on how these values could affect the task identifiers.

378 1/0/0 1/0/0 0/1/0 0/0/1 0/1/0 0/0/1 20 12 0.50/0.50/0.00 0.50/0.50/0.00 10 15 0.50/0.00/0.50 0.50/0.00/0.50 0.00/0.50/0.50 0 00/0 50/0 50 .33/0.33/0.33 0.33/0.33/0.33 10 382 lent 5 -10 -11 -20 -10 10 20 -10 (a) copy (op1) / copy (op2) / op1+op2 (b) to\_fr / to\_de / to\_it

Figure 4: Task vectors of Llama-3 8B projected onto two axes chosen by LDA for two sets of tasks: (a) copy (op1), copy (op2) and op1+op2 and (b) to\_fr, to\_de and to\_it. For tasks  $t_1, t_2, t_3$ , we use " $\mathbb{P}(t_1)/\mathbb{P}(t_2)/\mathbb{P}(t_3)$ " to denote different levels of task mixtures, e.g., "0.50/0.50/0.00" represents the case where the in-context task examples are  $50\% t_1$ ,  $50\% t_2$  and  $0\% t_3$ .

For each set of tasks, we collect the task vectors for each individual task and task vectors extracted from prompts that contain examples of different tasks. In Figure 4, we project task vectors along two axes chosen by linear discriminant analysis (LDA).

**Finding 3:** *LLMs internally combine task vectors during task superposition.* 

Interestingly, we observe that the locations of task vectors of a mixture of tasks strongly correlate with the locations of task vectors for each individual task and the in-context task example distribution (the mixture ratio for examples of different tasks). For example, if the prompt includes an equal number of in-context examples from each task, the task vectors are roughly centered in the middle; if the prompt only contains in-context examples of two tasks, then the task vectors roughly lie on the connecting line between task vectors of two individual tasks. We argue that this observation is indicative of the fact that, when prompted with a mixture of in-context task examples, LLMs internally combine task vectors.

As we observe signs that LLMs internally compose task vectors, we want to further investigate whether 411 we can reproduce the task superposition phenomenon by patching in a convex combination of task 412 vectors. For example, for tasks copy (op1) and copy (op2), we first extract the corresponding 413 task vectors  $V_{\text{copy (op1)}}$  and  $V_{\text{copy (op2)}}$  on Llama-3 8B using the method described in Appendix C. 414 We then make a convex combination of the two task vectors with parameter  $\lambda$  that controls the ratio: 415

$$\mathbf{V}_{\text{interpolate},\lambda} = \lambda \cdot \mathbf{V}_{\text{copy(op1)}} + (1 - \lambda) \cdot \mathbf{V}_{\text{copy(op2)}}$$

**Finding 4:** Convex combinations of task vectors produce task superposition.

420 For a new query (in this scenario in the form "{op1}@{op2}="), we patch the vector V<sub>interpolate, $\lambda$ </sub> 421 into the model at the task vector layer. We calculate the model output probabilities that correspond to 422 each task while we vary  $\lambda$ . For each  $\lambda$ , we use 100 different queries and plot the average probabilities 423 in the top row of Figure 5. As a comparison, in the bottom rows of Figure 5, we plot the corresponding 424 output probabilities when providing the models with prompts containing mixture of task examples 425 where the mixture ratio is controlled by  $\lambda$ .

In top row of Figure 5, we observe that patching convex combinations of task vectors into the model 427 produces task superposition. We would also like to point out that in Figure 5b, although irrelevant 428 outputs sum up to a large probability, the task answers for two tasks to\_de and to\_it in most cases 429 will still be the top-2 answers. 430

Comparing the top rows and the bottom rows, we can see that top rows (the scenario of interpolating 431 task vectors of individual tasks) have larger probabilities of irrelevant output (category other).



379

380

381

384 385 386

387 388

389 390

391 392

393

394

396

397

398

399 400

401 402

403

404

405

406

407

408

409

410

417 418 419



(a) Tasks: copy (op1) and copy (op2)



Figure 5: On Llama-3 8B, we vary the proportion,  $\lambda$ , between two tasks and observe how the output probabilities for the correct answers change. The proportion  $\lambda$  is varied in two ways: (1) in the top row, we plot the output from patching in a convex combination of task vectors for two tasks. (2) in the bottom row, we plot the output from a mixed proportion of in-context examples for the two tasks. Subplot (a) shows the output probabilities from mixing two copy tasks and (b) shows the probabilities from mixing two translate tasks.

Task vector interpolation also produces less of a linear relationship between  $\lambda$  and the output probabilities. This shows that while convex combinations of task vectors are sufficient for producing task superposition, this does not fully explain task superposition. We leave it to future work to investigate other mechanistic explanations of task superposition.

# 7 TASK SUPERPOSITION CAPABILITIES AS THE MODEL SCALES

**Finding 5:** Within the same LLM family, bigger models can solve more tasks in parallel and better calibrate to ICL distribution.

We want to further investigate how models' task superposition capabilities changes as the model size scales. In particular, we investigate two questions: 1) whether larger models can perform more tasks in-context and 2) whether larger models can align their output distribution more closely with the distribution of task examples provided in the prompt. We chose the Qwen-1.5 model family since it contains several model sizes ranging from 0.5B to 14B parameters.

We first introduce a quantity which captures the capability of a model to perform multiple tasks. Given a prompt that contains examples of K tasks, we define r to be the number of these tasks whose correct answers appear among the model's top-K most likely outputs. Note that  $r \le K$ .

To see how close the model align the output distribution with the distribution of task examples, we use KL-divergence defined below:

$$\mathrm{KL}(\mathcal{P}||\mathcal{D}) = \sum_{x \in X} \mathcal{P}(x) \log\left(\frac{\mathcal{P}(x)}{\mathcal{D}(x)}\right),\tag{2}$$

where  $\mathcal{P}$  is the models' probabilities on the outputs when correctly performing each task on the query and  $\mathcal{D}$  is the in-context task example distribution. For example the prompt in Figure 1a (left) gives  $\mathcal{P} = [0.5217, 0.1316, 0.1110, 0.2169, ...]$  and  $\mathcal{D} = [0.25, 0.25, 0.25, 0.25, 0, ...].$ 

We consider the setting of K = 6 different tasks: given an input of the form "{num} $\rightarrow$ " where num is a two-digit integer, we consider 6 tasks that output (1) num itself, (2) negation of num, (3) num + 1, (4) num - 1, (5) num  $\times$  2 and (6) num<sup>2</sup>.

486 We choose the number of in-context examples m = 60 (each task has 10 examples) and configure the prompt with three different in-context task example distributions  $\mathcal{D}_1, \mathcal{D}_2$  and  $\mathcal{D}_3$ . In particular, 488  $\mathcal{D}_1$  is the uniform distribution,  $\mathcal{D}_2$  has probability 0.5 on the third task and 0.1 on other tasks, and 489  $\mathcal{D}_3$  is a distribution with probabilities alternating between 0.25 and 0.083.

For each in-context task examples distribution  $\mathcal{D}_i$ , we generate 100 prompts and for each prompt we calculate the probabilities of outputs when correctly performing each task. The average values of rand KL-divergence under three distributions are shown in Figure 6.



Figure 6: (a) Average number of tasks completed, r, and (b) KL divergence for Qwen-1.5 model family 504 under ICL distributions  $\mathcal{D}_1, \mathcal{D}_2$  and  $\mathcal{D}_3$  where  $\mathcal{D}_1$  is the uniform distribution,  $\mathcal{D}_2$  has probability 0.5 505 on the third task and 0.1 on other tasks, and  $\mathcal{D}_3$  is a distribution with probabilities alternating between 506 0.25 and 0.083.

507 In Figure 6a, we can observe that bigger models have higher r values (except for task distribution 508  $\mathcal{D}_2$ , 4B model has slightly lower r than that of the 1.8B model). This shows bigger models will 509 have more correct answers of tasks show up in their top-K probable outputs and therefore they can 510 solve more tasks at the same time. In Figure 6b, we can see that for larger models like Qwen-1.5 511 7B and Qwen-1.5 14B, the KL-divergence values are small, and for each model, the differences 512 between KL-divergence values under in-context task example distributions  $\mathcal{D}_1$ ,  $\mathcal{D}_2$  and  $\mathcal{D}_3$  are small. 513 This indicates that bigger models can better calibrate their output distribution to the in-context task 514 example distribution.

515 516

517

487

490

491

492

493 494

495

496

497

498

499

500 501

502

#### 8 LIMITATIONS AND FUTURE DIRECTIONS

518 One limitation of our work is the current gap between the demonstrated capability of LLMs to 519 perform task superposition and its practical application in real-world scenarios. While we have shown 520 that LLMs possess the capacity to execute multiple tasks simultaneously, conventional decoding 521 algorithms are not equipped to fully leverage this capability. This limitation stems from what we 522 term "generation collapse," a phenomenon where, after the first token is generated, the model tends 523 to converge on predicting tokens for a single task, effectively negating its ability for multi-task execution. 524

525 This collapse presents a substantial challenge in harnessing the full power of task superposition. It 526 highlights a critical area for future research: developing decoding strategies that can maintain the 527 model's multi-task state throughout the generation process. Recent work by Shen et al. (2024b) offers 528 some hope that this direction may be fruitful, by proposing a "superposed decoding" algorithm. Their 529 method efficiently generates multiple streams of tokens from a single inference pass by utilizing superposed token embeddings. While this approach represents a significant step forward, it also 530 highlights the potential for further innovation in this area. 531

532 533

534

#### 9 CONCLUSION

We report on the discovery of task superposition, which is the ability of LLMs to simultaneously solve distinct tasks from in-context examples. Task superposition is present in a variety of pretrained 537 models, and becomes more accurate at predicting the distribution of tasks as the model size increases. We also find evidence that while displaying task superposition, models internally mix the task vectors 538 of each individual task. We hope that our findings will contribute to understanding in-context learning mechanisms and enhance our knowledge of LLMs overall.

#### 540 REFERENCES 541

547

567

581

582

583

584

585

586

589

Rishabh Agarwal, Avi Singh, Lei M. Zhang, Bernd Bohnet, Stephanie Chan, Ankesh Anand, Zaheer 542 Abbas, Azade Nova, John D. Co-Reyes, Eric Chu, Feryal Behbahani, Aleksandra Faust, and Hugo 543 Larochelle. Many-shot in-context learning, 2024. 544

- Kwangjun Ahn, Xiang Cheng, Hadi Daneshmand, and Suvrit Sra. Transformers learn to implement 546 preconditioned gradient descent for in-context learning. arXiv preprint arXiv:2306.00297, 2023.
- 548 AI@Meta. Llama 3 model card. 2024. URL https://github.com/meta-llama/llama3/ blob/main/MODEL\_CARD.md. 549
- 550 Ekin Akyürek, Dale Schuurmans, Jacob Andreas, Tengyu Ma, and Denny Zhou. What learning 551 algorithm is in-context learning? Investigations with linear models, May 2023. URL http: 552 //arxiv.org/abs/2211.15661. arXiv:2211.15661 [cs]. 553
- 554 Jinze Bai, Shuai Bai, Yunfei Chu, Zeyu Cui, Kai Dang, Xiaodong Deng, Yang Fan, Wenbin Ge, Yu Han, Fei Huang, Binyuan Hui, Luo Ji, Mei Li, Junyang Lin, Runji Lin, Dayiheng Liu, Gao Liu, 555 Chengqiang Lu, Keming Lu, Jianxin Ma, Rui Men, Xingzhang Ren, Xuancheng Ren, Chuanqi Tan, 556 Sinan Tan, Jianhong Tu, Peng Wang, Shijie Wang, Wei Wang, Shengguang Wu, Benfeng Xu, Jin Xu, An Yang, Hao Yang, Jian Yang, Shusheng Yang, Yang Yao, Bowen Yu, Hongyi Yuan, Zheng 558 Yuan, Jianwei Zhang, Xingxuan Zhang, Yichang Zhang, Zhenru Zhang, Chang Zhou, Jingren Zhou, 559 Xiaohuan Zhou, and Tianhang Zhu. Qwen technical report. arXiv preprint arXiv:2309.16609, 2023a. 561
- 562 Yu Bai, Fan Chen, Huan Wang, Caiming Xiong, and Song Mei. Transformers as statisticians: 563 Provable in-context learning with in-context algorithm selection, 2023b.
- 564 Yu Bai, Heyan Huang, Cesare Spinoso-Di Piano, Marc-Antoine Rondeau, Sanxing Chen, Yang Gao, 565 and Jackie Chi Kit Cheung. Identifying and analyzing task-encoding tokens in large language 566 models. (arXiv:2401.11323), February 2024. doi: 10.48550/arXiv.2401.11323. URL http: //arxiv.org/abs/2401.11323. arXiv:2401.11323 [cs]. 568
- Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhari-569 wal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, 570 Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel 571 Ziegler, Jeffrey Wu, Clemens Winter, Chris Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, 572 Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Rad-573 ford, Ilya Sutskever, and Dario Amodei. Language Models are Few-Shot Learners. In Ad-574 vances in Neural Information Processing Systems, volume 33, pp. 1877–1901. Curran Asso-575 ciates, Inc., 2020. URL https://proceedings.neurips.cc/paper\_files/paper/ 576 2020/hash/1457c0d6bfcb4967418bfb8ac142f64a-Abstract.html. 577
- Stephanie Chan, Adam Santoro, Andrew Lampinen, Jane Wang, Aaditya Singh, Pierre Richemond, 578 James McClelland, and Felix Hill. Data distributional properties drive emergent in-context learning 579 in transformers. Advances in Neural Information Processing Systems, 35:18878–18891, 2022. 580
  - Brian Cheung, Alex Terekhov, Yubei Chen, Pulkit Agrawal, and Bruno Olshausen. Superposition of many models into one, June 2019. URL http://arxiv.org/abs/1902.05522. arXiv:1902.05522 [cs].
  - Hyeong Kyu Choi and Yixuan Li. Picle: Eliciting diverse behaviors from large language models with persona in-context learning. (arXiv:2405.02501), May 2024. doi: 10.48550/arXiv.2405.02501. URL http://arxiv.org/abs/2405.02501. arXiv:2405.02501 [cs].
- 588 Damai Dai, Yutao Sun, Li Dong, Yaru Hao, Zhifang Sui, and Furu Wei. Why can gpt learn incontext? language models secretly perform gradient descent as meta optimizers. arXiv preprint 590 arXiv:2212.10559, 2022.
- Tuan Dinh, Yuchen Zeng, Ruisu Zhang, Ziqian Lin, Shashank Rajput, Michael Gira, Jy-yong 592 Sohn, Dimitris Papailiopoulos, and Kangwook Lee. Lift: Language-interfaced fine-tuning for non-language machine learning tasks. arXiv preprint arXiv:2206.06565, 2022.

594 595 596 597 598	Nelson Elhage, Tristan Hume, Catherine Olsson, Nicholas Schiefer, Tom Henighan, Shauna Kravec, Zac Hatfield-Dodds, Robert Lasenby, Dawn Drain, Carol Chen, Roger Grosse, Sam McCandlish, Jared Kaplan, Dario Amodei, Martin Wattenberg, and Christopher Olah. Toy Models of Superpo- sition, September 2022. URL http://arxiv.org/abs/2209.10652. arXiv:2209.10652 [cs].
599 600 601 602	Fabian Falck, Ziyu Wang, and Chris Holmes. Is in-context learning in large language models bayesian? a martingale perspective. (arXiv:2406.00793), June 2024. URL http://arxiv.org/abs/2406.00793. arXiv:2406.00793 [cs, stat].
603 604 605	Angeliki Giannou, Shashank Rajput, Jy-yong Sohn, Kangwook Lee, Jason D. Lee, and Dimitris Papailiopoulos. Looped transformers as programmable computers, 2023.
606 607 608	Riccardo Grazzi, Julien Siems, Simon Schrodi, Thomas Brox, and Frank Hutter. Is mamba capable of in-context learning? (arXiv:2402.03170), April 2024. doi: 10.48550/arXiv.2402.03170. URL http://arxiv.org/abs/2402.03170. arXiv:2402.03170 [cs].
609 610 611	Albert Gu and Tri Dao. Mamba: Linear-time sequence modeling with selective state spaces. (arXiv:2312.00752), December 2023. doi: 10.48550/arXiv.2312.00752. URL http://arxiv. org/abs/2312.00752. arXiv:2312.00752 [cs].
613 614 615	Roee Hendel, Mor Geva, and Amir Globerson. In-Context Learning Creates Task Vectors, October 2023. URL http://arxiv.org/abs/2310.15916. arXiv:2310.15916 [cs].
616 617 618	Gabriel Ilharco, Marco Tulio Ribeiro, Mitchell Wortsman, Suchin Gururangan, Ludwig Schmidt, Hannaneh Hajishirzi, and Ali Farhadi. Editing Models with Task Arithmetic, March 2023. URL http://arxiv.org/abs/2212.04089. arXiv:2212.04089 [cs].
619 620 621	Janus. Simulators, 2022. URL https://www.lesswrong.com/posts/ vJFdjigzmcXMhNTsx/.
622 623 624	Andrew K Lampinen, Ishita Dasgupta, Stephanie CY Chan, Kory Matthewson, Michael Henry Tessler, Antonia Creswell, James L McClelland, Jane X Wang, and Felix Hill. Can language models learn from explanations in context? <i>arXiv preprint arXiv:2204.02329</i> , 2022.
625 626 627 628	Michael Laskin, Luyu Wang, Junhyuk Oh, Emilio Parisotto, Stephen Spencer, Richie Steigerwald, DJ Strouse, Steven Hansen, Angelos Filos, Ethan Brooks, et al. In-context reinforcement learning with algorithm distillation. <i>arXiv preprint arXiv:2210.14215</i> , 2022.
629 630 631	Tianle Li, Ge Zhang, Quy Duc Do, Xiang Yue, and Wenhu Chen. Long-context llms struggle with long in-context learning. (arXiv:2404.02060), April 2024. doi: 10.48550/arXiv.2404.02060. URL http://arxiv.org/abs/2404.02060. arXiv:2404.02060 [cs].
632 633 634	Ziqian Lin and Kangwook Lee. Dual Operating Modes of In-Context Learning, February 2024. URL http://arxiv.org/abs/2402.18819. arXiv:2402.18819 [cs].
635 636 637 638	Sheng Liu, Haotian Ye, Lei Xing, and James Zou. In-context vectors: Making in context learning more effective and controllable through latent space steering. (arXiv:2311.06668), February 2024. doi: 10.48550/arXiv.2311.06668. URL http://arxiv.org/abs/2311.06668. arXiv:2311.06668 [cs].
639 640 641	Xinxi Lyu, Sewon Min, Iz Beltagy, Luke Zettlemoyer, and Hannaneh Hajishirzi. Z-icl: Zero-shot in-context learning with pseudo-demonstrations. <i>arXiv preprint arXiv:2212.09865</i> , 2022.
642 643 644 645 646	Sewon Min, Xinxi Lyu, Ari Holtzman, Mikel Artetxe, Mike Lewis, Hannaneh Hajishirzi, and Luke Zettlemoyer. Rethinking the role of demonstrations: What makes in-context learning work? (arXiv:2202.12837), October 2022. doi: 10.48550/arXiv.2202.12837. URL http://arxiv.org/abs/2202.12837. arXiv:2202.12837 [cs].
0.47	moire Language models are multiverse generators, January 2021 JIRL https://gonorativo

647 moire. Language models are multiverse generators, January 2021. URL https://generative. ink/posts/language-models-are-multiverse-generators/.

648 640	Vishvak Murahari, Carlos E. Jimenez, Runzhe Yang, and Karthik Narasimhan. Datamux: Data
650	multiplexing for neural networks. (arXiv:2202.09318), November 2022. doi: 10.48550/arXiv.2202. 09318. URL http://arxiv.org/abs/2202.09318. arXiv:2202.09318 [cs].
651	
652	Cleo Nardo. The waluigi effect (mega-post), 2023. URL https://www.lesswrong.com/
653	posts/D/PumeYTDPiBTp31//the-waluigi-effect-mega-post.
654	Catherine Olsson, Nelson Elhage, Neel Nanda, Nicholas Joseph, Nova DasSarma, Tom Henighan,
655	Ben Mann, Amanda Askell, Yuntao Bai, Anna Chen, et al. In-context learning and induction heads.
656	arXiv preprint arXiv:2209.11895, 2022.
657	Johannes Von Oswald, Ewind Niklasson, Ettore Pandazzo, Joao Sacramento, Alexander Mordvintsev
658	Andrey Zhmoginov and Max Vladymyrov Transformers Learn In-Context by Gradient Descent In
659	Proceedings of the 40th International Conference on Machine Learning, pp. 35151–35174. PMLR,
661	July 2023. URL https://proceedings.mlr.press/v202/von-oswald23a.html.
660	ISSN: 2640-3498.
662	Madhur Danwar, Kabir Abuja, and Navin Goval. In contact learning through the bayesian prism
664	June 2023. URL https://arxiv.org/abs/2306.04891v2.
665	
666	Jongno Park, Jaeseung Park, Zneyang Xiong, Nayoung Lee, Jaewoong Cho, Samet Oymak, Kang-
667	in-context learning tasks (arXiv:2402.04248) April 2024 doi: 10.48550/arXiv.2402.04248 URI
668	http://arxiv.org/abs/2402.04248. arXiv:2402.04248 [cs].
669	
670	Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, Ilya Sutskever, et al. Language
671	models are unsupervised multitask learners. OpenAI blog, 1(8):9, 2019.
672	Allan Raventos, Mansheej Paul, Feng Chen, and Surya Ganguli. The effects of pretraining task
673	diversity on in-context learning of ridge regression. In ICLR Workshop on Mathematical and
674	Empirical Understanding of Foundation Models (ME-FoMo), 2023.
675	Gautam Reddy. The mechanistic basis of data dependence and abrupt learning in an in-context
677	classification task. In <i>The Twelfth International Conference on Learning Representations</i> , 2024.
679	URL https://openreview.net/forum?id=aN4Jf6Cx69.
679	
680	Laria Reynolds and Kyle McDonell. Multiversal views on language models, February 2021. URL
681	http://arxiv.org/abs/2102.00591. arXiv.2102.00591 [cs].
682	Murray Shanahan, Kyle McDonell, and Laria Reynolds. Role play with large language models.
683	Nature, 623(7987):493–498, 2023.
684	Ethan Shen, Alan Fan, Sarah M. Pratt Jae Sung Park Matthew Wallingford Sham M Kakade Ari
685	Holtzman, Ranjav Krishna, Ali Farhadi, and Aditva Kusupati. Superposed decoding: Multiple
686	generations from a single autoregressive inference pass. (arXiv:2405.18400), May 2024a. doi: 10.
687	48550/arXiv.2405.18400. URL http://arxiv.org/abs/2405.18400. arXiv:2405.18400
688	[cs].
689	Ethan Shen, Alan Fan, Sarah M Pratt, Jae Sung Park, Matthew Wallingford, Sham M Kakade, Ari
690	Holtzman, Ranjay Krishna, Ali Farhadi, and Aditya Kusupati. Superposed decoding: Multiple
691	generations from a single autoregressive inference pass. arXiv preprint arXiv:2405.18400, 2024b.
602	Erio Todd Millicont I. Li Arnoh Con Chorme, Acres Marillan Davier C. Wallace and D. (1D)
694	Enc roud, Minicent L. Li, Arnao Sen Sharma, Aaron Muener, Byron C. Wallace, and David Bau. Function vectors in large language models (arXiv:2310.15213) February 2024. doi: 10.48550/
695	arXiv.2310.15213. URL http://arxiv.org/abs/2310.15213. arXiv:2310.15213 [cs]
696	
697	Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay
698	Basniykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, et al. Llama 2: Open foundation and fine tuned chat models, arViv preprint arViv:2207.00288, 2023
699	and mic-tuned enat models. arxiv preprint arxiv:2507.09268, 2025.
700	Colin Wei, Yining Chen, and Tengyu Ma. Statistically meaningful approximation: a case study on
701	approximating turing machines with transformers. <i>Advances in Neural Information Processing Systems</i> , 35:12071–12083, 2022.

702 703 704 705	Jerry Wei, Jason Wei, Yi Tay, Dustin Tran, Albert Webson, Yifeng Lu, Xinyun Chen, Hanxiao Liu, Da Huang, Denny Zhou, et al. Larger language models do in-context learning differently. <i>arXiv</i> preprint arXiv:2303.03846, 2023.
705 706 707 708	Jingfeng Wu, Difan Zou, Zixiang Chen, Vladimir Braverman, Quanquan Gu, and Peter L Bartlett. How many pretraining tasks are needed for in-context learning of linear regression? In <i>International</i> <i>Conference on Learning Representations (ICLR)</i> , 2024.
709 710 711	Sang Michael Xie, Aditi Raghunathan, Percy Liang, and Tengyu Ma. An explanation of in-context learning as implicit bayesian inference. (arXiv:2111.02080), July 2022. doi: 10.48550/arXiv.2111. 02080. URL http://arxiv.org/abs/2111.02080. arXiv:2111.02080 [cs].
712 713 714	Liu Yang, Kangwook Lee, Robert Nowak, and Dimitris Papailiopoulos. Looped transformers are better at learning learning algorithms. <i>arXiv preprint arXiv:2311.12424</i> , 2023.
715 716 717	Yuchen Zeng, Wonjun Kang, Yicong Chen, Hyung Il Koo, and Kangwook Lee. Can mllms perform text-to-image in-context learning? (arXiv:2402.01293), April 2024. doi: 10.48550/arXiv.2402. 01293. URL http://arxiv.org/abs/2402.01293. arXiv:2402.01293 [cs].
718 719 720 721 722	Yufeng Zhang, Fengzhuo Zhang, Zhuoran Yang, and Zhaoran Wang. What and how does in- context learning learn? bayesian model averaging, parameterization, and generalization. (arXiv:2305.19420), October 2023. doi: 10.48550/arXiv.2305.19420. URL http://arxiv. org/abs/2305.19420. arXiv:2305.19420 [cs, stat].
723 724	Hattie Zhou, Azade Nova, Hugo Larochelle, Aaron Courville, Behnam Neyshabur, and Hanie Sedghi. Teaching algorithmic reasoning via in-context learning. <i>arXiv preprint arXiv:2211.09066</i> , 2022.
725 726 727 728	Hattie Zhou, Arwen Bradley, Etai Littwin, Noam Razin, Omid Saremi, Josh Susskind, Samy Bengio, and Preetum Nakkiran. What algorithms can transformers learn? a study in length generalization. <i>arXiv preprint arXiv:2310.16028</i> , 2023.
729	
730	
732	
733	
734	
735	
736	
737	
738	
739	
740	
741	
742	
743	
744	
745	
746	
747	
748	
749	
750	
751	
752	
103	
755	
100	

# 756 A NOTATIONS 757

758		
759	Notation	Description
760	K	Number of tasks
761	l	Length of a task's output
762	$\ell$	layer $\ell$ for a model
763	m	Number of in-context examples
763	n	Length of each in-context example
764	$\nu$	loken vocabulary
705	$g_{i(\cdot)}$	Operation performed by Task i
766	$x^{(j)}_{(i)}$	Data for example $j$
767	$oldsymbol{y}^{(j)}$	Label for example $j$
768	$s_m$	m in-context examples
769	$oldsymbol{f}(\cdot)$	Model (predictor)
770	<i>p</i>	Positional encodings
771		
772		
773		
774		
775		
776		
777		
778		
770		
790		
701		
781		
782		
783		
784		
785		
786		
787		
788		
789		
790		
791		
792		
793		
794		
795		
796		
797		
709		
790		
799		
800		
801		
802		
803		
804		
805		
806		
807		
808		
809		

# 810 B IMPLEMENTATION DETAILS ON CALCULATING PROBABILITIES

812 In this section we provide details on how we calculate probabilities of different outputs given a prompt in our setting.
 814

**Notations.** Let  $\mathcal{V}$  be the token vocabulary, LM be an LLM, T be the tokenizer. We use "..." to represent a string, <...> to represent a single token where the content within the angle brackets is an integer representing token's index in vocabulary. For example, token <266> corresponds to "at". We use [<...>, ..., <...>] to represent a sequence of tokens. Given a tokenizer, we use two functions tok(·) and detok(·) to tokenize strings and detokenize tokens. For example tok("superposition") = [<9712>,<3571>] and detok([<16>,<10>,<16>,<28>,<17>]) = "1+1=2".

In our in-context learning setting, an input string consists of in-context examples (separated by the delimiter "\n") and a query. For example, an example prompt can be " $1+1=2\n2+2=4\n3+3=$ ".

We view an LLM as a next-token predictor that outputs a probability distribution over the token space given input and there is a corresponding  $\mathbb{P}(\cdot|\cdot)$  such that given a sequence of tokens  $[v_1, ..., v_M]$ where  $v_j \in \mathcal{V}$ ,  $\mathbb{P}(u \mid [v_1, ..., v_M])$  measures the probability of the next token being u where  $u \in \mathcal{V}$ .

828 Measuring the probabilities of task answers. Let *I* be the input prompt. For example, in the 829 example in Figure 1a (left), the prompt is " $11+26->37\n33+13->quarante-six\n ...30+25->fifty 830 five\n91+83->". We consider four tasks: 1) numerical addition, 2) addition in English, 3) addition$ 831 in French and 4) addition in Spanish. The corresponding task answers (the output of correctly832 performing task on the query) are "174", "one hundred and seventy-four", "cent soixante-quatorze"833 and "ciento setenta y cuatro", respectively. We want to measure the probability of each task answer.

Let o be a task answer in string. Let  $[v_1, ..., v_M] := tok(I)$  and let  $[u_1, ..., u_N] := tok(o)$ . Then the probability of the task answer o given prompt I can be calculated as

16

$$\mathbb{P}(u_1 \mid [v_1, ..., v_M]) \prod_{j=2}^N \mathbb{P}(u_j \mid [v_1, ..., v_M, u_1, ..., u_{j-1}]).$$
(3)

837 838

839 840

841

842 843

# 864 C IMPLEMENTATION DETAILS ON TASK VECTORS

We use the task vector definition from Hendel et al. (2023). For example, for task copy (op1) in Figure 1b (left), the procedure to collect the task vector consists of

- 1. Collect a dataset of 100 ICL sample prompts. Each prompt consists of m = 60 in-context examples of a particular task and a query  $\boldsymbol{x}^{(m+1)}$ . Each task example  $(\boldsymbol{x}^{(j)}, \boldsymbol{y}^{(j)})$  follows the form "{op1}@{op2}={op1}", where  $\boldsymbol{x}^{(j)}$  has the form "{op1}@{op2}=" and  $\boldsymbol{y}^{(j)}$  is performing task copy (op1) on  $\boldsymbol{x}^{(j)}$ , namely op1.
- For each prompt ≡ s = [x<sup>(1)</sup>, y<sup>(1)</sup>, ..., x<sup>(m)</sup>, y<sup>(m)</sup>, x<sup>(m+1)</sup>] in the dataset, we feed s into the transformer model f, and extract the feature (which is a vector) at the last "=" token in layer l. Call this vector f(s; l). Then we average f(s; l) across all prompt s to get v(l) for layer l.
- 3. Now for each layer ℓ we have a vector v(ℓ). We run a forward pass with one query x in the form "{op1}@{op2}=" and we patch in v(ℓ) at the "=" token position in layer ℓ, simulating the effect of a complete context. We repeat this process 100 times for different query x and get an accuracy acc<sub>ℓ</sub> of performing task copy (op1) with vector v(ℓ).
  - 4. The task vector layer  $\ell^*$  is selected by

$$\ell^* = \arg\max_{\ell} acc_{\ell},$$

and we define the task vector  $V_{\text{copy(op1)}} := v = v(\ell^*)$ .

Here we record the task vector layer where task vectors are extracted in Section 6.

Task	Task vector layer
<pre>copy(op1), copy(op2), op1+op2</pre>	14
<pre>to_de(op1),to_fr(op1),to_it(op1)</pre>	19

Table 1: Task vector layer for various tasks considered in Section 6.

#### 918 D CONSTRUCTION DISPLAYING SUPERPOSITION 919

920 In this section we construct a Transformer that is performing superposition of multiple tasks at 921 inference. For this purpose, we first construct a Transformer that copies from *n*-tuple in-context 922 examples the *i*-th one, as well as any function using the ReLU layers. We then create indicator 923 vectors, for each task, which show whether a specific task is present in-context or not. As a last step, 924 we combine these indicator vectors to create the superposition of different tasks. Notice that using the parallel heads of the transformer architecture we can process each task independently until the 925 926 last step in which the predictions are combined.

# D.1 OVERVIEW

927

928

931 932

944

947

929 Here we provide a brief overview of how the construction is implemented, while latter we provide 930 the corresponding details.

**Prediction based on multiple tasks.** Assume that we are given m in-context samples 933  $(\boldsymbol{x}_1^{(j)},\ldots,\boldsymbol{x}_{n-2}^{(j)}, \textbf{`='}, \boldsymbol{y}^{(j)})_{j=1}^m$  where `=' represents a specific value used only for preceding the label, and a set of k different Transformers TF<sub>i</sub> which can implement the T different desired tasks, 934 935 where each deterministic task is denoted as  $g_i(x^{(j)})$  with  $i \in [k]$  and  $j \in [m]$ , *i.e.*  $y^{(j)} = g_i(x^{(j)})$ 936 for some task i dependent on sample j. Using the weights of each  $TF_i$ , we can compute outputs of 937 the following form: 938

939	ſ	$x_1^{(j)}$	 $x_{n-2}^{(j)}$	=	$oldsymbol{y}^{(j)}$	]	<b>[</b>	$m{x}_1^{(j)}$	 $oldsymbol{x}_{n-2}^{(j)}$	=	$oldsymbol{y}^{(j)}$	]
940 941		Ō	 0	0	0	_		0	 0	0	$\ m{g}_1(m{x}^{(j)}) - m{y}^{(j)}\ _1$	
942		÷	÷	÷	÷			÷	÷	÷	÷	
943	[	0	 0	0	0	]	L	0	 0	0	$\ m{g}_T(m{x}^{(j)}) - m{y}^{(j)}\ _1$	· · · ]

We use the  $l_1$  norm to aggregate the prediction, in case that the task is multi-dimensional. These 945 differences are used to identify tasks, as  $\|\boldsymbol{g}_i(\boldsymbol{x}^{(j)}) - \boldsymbol{y}^{(j)}\|_1 \approx 0$  for  $\boldsymbol{y}^{(j)}$  coming from task *i*. Different 946 heads at each layer in the model are used to execute each of the tasks in parallel using the weights from TF<sub>i</sub>. In Appendix D we construct tasks where an arbitrary function  $g_i(x_i^{(j)})$  is implemented 948 using ReLUs for some fixed *l* that is task-specific. 949

950 Creating task identifiers. Having the differences between the implemented function and the 951 label, we first use the ReLUs to clean up the vectors  $v_k$  so that only the positions in each vector 952 that are associated with a task are maintained and the rest are set to  $1^4$ . We thus create the vectors 953  $(v'_k)_i = \|g_k(x_{1:l-1}^{(j)}) - x_l^{(j)}\|_1$  and  $(v'_*)_* = 1$  otherwise. Now we use ReLUs to threshold and create 954 an indicator vectors  $\mathbf{1}_{\{\|\boldsymbol{g}_k(\boldsymbol{x}_{1:l-1}^{(j)})-\boldsymbol{x}_l^{(j)}\|_1\approx 0\}}$  which identify the task, *i.e.* these are task identifiers. 955 Notice that if the task is correctly predicted then the difference should be close to 0 (up to some 956 error), while if the task is not identified the corresponding value would not be 0; the rest of the 957 rows would be 1. We have created one vector for each task, which has 1 in the position of the 958 corresponding task if the task was identified in the context. 959

960 Averaging and task superposition. As a last step, we average all the task identifiers and place the 961 result in the last column, in which the next prediction will happen. We then use the averaged task 962 identifier to weight the prediction of each task based on it, as in task superposition. If the task has 963 been identified multiple times in the context, it would be assigned a higher weight/probability.

965 **D.2** TASK IDENTIFICATION

The first task for performing task superposition based on in-context examples is to define a set of 967 tasks that the model is able to implement. 968

- 969 First, the outputs of tasks need to be identified.
- 970 971

964

<sup>&</sup>lt;sup>4</sup>This step is not mandatory, but it ensures that we have no values over which we have no control. We leave as future work an error analysis on how these values could affect the task identifiers

**Lemma 1.** Consider the following input

  $\boldsymbol{X} = \begin{bmatrix} \boldsymbol{x}_1^{(1)} & \dots & \boldsymbol{y}^{(j-1)} & \boldsymbol{x}_1^{(j)} & \dots & \boldsymbol{x}_{n-2}^{(j)} & = & \boldsymbol{y}^{(j)} & \boldsymbol{x}_1^{(j+1)} & \dots \\ 0 & \dots & 0 & 0 & \dots & 0 & 1 & 0 & 0 & \dots \\ \boldsymbol{0} & \dots & \boldsymbol{0} & \boldsymbol{0} & \dots & \boldsymbol{0} & \boldsymbol{0} & \boldsymbol{0} & \dots \end{bmatrix} ,$ 

where  $x_i^{(j)} \in \mathbb{R}^{d-1}$  before the positional encodings are added, with one additional dimension that represents if the symbol is an 'equals' symbol. Then, a 1-layer transformer with a single attention head and embedding dimension  $\mathcal{O}(d + \log(mn))$  can output

$$\boldsymbol{X} = \begin{bmatrix} \boldsymbol{x}_1^{(1)} & \dots & \boldsymbol{y}^{(j-1)} & \boldsymbol{x}_1^{(j)} & \dots & \boldsymbol{x}_{n-2}^{(j)} & = & \boldsymbol{y}^{(j)} & \boldsymbol{x}_1^{(j+1)} & \dots \\ 0 & \dots & 1 & 0 & \dots & 0 & 0 & 1 & 0 & \dots \end{bmatrix}$$

*Proof.* With positional encodings appended, let the input have the following structure:

$$\boldsymbol{X} = \begin{bmatrix} \boldsymbol{x}_{1}^{(1)} & \dots & \boldsymbol{x}_{1}^{(j)} & \boldsymbol{x}_{2}^{(j)} & \dots & \boldsymbol{x}_{n-2}^{(j)} & = & \boldsymbol{y}^{(j)} & \boldsymbol{x}_{1}^{(j+1)} & \dots \\ 0 & \dots & 0 & 0 & \dots & 0 & 1 & 0 & 0 & \dots \\ \boldsymbol{p}_{n+1} & \dots & \boldsymbol{p}_{jn+1} & \boldsymbol{p}_{jn+2} & \dots & \boldsymbol{p}_{jn+n-2} & \boldsymbol{p}_{jn+n-1} & \boldsymbol{p}_{jn+n} & \boldsymbol{p}_{(j+1)n+1} & \dots \\ \boldsymbol{p}_{n} & \dots & \boldsymbol{p}_{jn} & \boldsymbol{p}_{jn+1} & \dots & \boldsymbol{p}_{jn+n-3} & \boldsymbol{p}_{jn+n-2} & \boldsymbol{p}_{jn+n-1} & \boldsymbol{p}_{(j+1)n} & \dots \end{bmatrix}$$
(4)

To rotate the second row one position to the right, use the following matrices.

$$egin{aligned} m{W}_Q &= egin{bmatrix} m{0} & m{0} & m{0} & m{I} \end{bmatrix} \ m{W}_K &= egin{bmatrix} m{0} & m{0} & m{C} m{I} & m{0} \end{bmatrix} \ m{W}_V &= egin{bmatrix} m{0} & m{0} & m{0} & m{0} & m{0} \end{bmatrix} \ m{W}_V &= egin{bmatrix} m{0} & m{0} & m{0} & m{0} \end{bmatrix} \ m{0} & m{0} & m{0} & m{0} \end{bmatrix} \ m{0} & m{0} & m{0} & m{0} \end{bmatrix} \ m{0} & m{0} & m{0} & m{0} \end{bmatrix} \ m{0} & m{0} & m{0} & m{0} \end{bmatrix} \ m{0} & m{0} & m{0} \end{bmatrix} \ m{0} & m{0} & m{0} \end{bmatrix} \ m{0} & m{0} & m{0} \end{bmatrix} \ m{0} & m{0} & m{0} \end{bmatrix} \ m{0} & m{0} & m{0} \end{bmatrix} \ m{0} & m{0} \end{bmatrix} \ m{0} & m{0} & m{0} \end{bmatrix} \ m{0} & m{0} & m{0} \end{bmatrix} \ m{0} & m{0} \end{bmatrix} \ m{0} & m{0} \end{bmatrix} \ m{0} & m{0} \end{bmatrix} \ m{0} & m{0} \end{bmatrix} \ m{0} \ m{0} \end{bmatrix} \ m{0} \ m{0} \end{bmatrix} \ m{0} \ m{0} \end{bmatrix} \ m{0} \ m{0} \end{bmatrix} \ m{0} \ m{0} \end{bmatrix} \ m{0} \ m{0} \end{bmatrix} \ m{0} \ m{0} \end{bmatrix} \ m{0} \ m{0} \end{bmatrix} \ m{0} \ m{0} \end{bmatrix} \ m{0} \ m{0} \end{bmatrix} \ m{0} \ m{0} \end{bmatrix} \ m{0} \ m{0} \end{bmatrix} \ m{0} \ m{0} \end{bmatrix} \ m{0} \ m{0} \end{bmatrix} \ m{0} \ m{0} \ m{0} \end{bmatrix} \ m{0} \ m{0} \end{bmatrix} \ m{0} \ m{0} \ m{0} \end{bmatrix} \ m{0} \ m{0} \ m{0} \ m{0} \end{bmatrix} \ m{0} \ m{$$

The pair  $W_Q$  and  $W_K$  attend tokens to the token directly to the right. The value matrix simply filters only the second row in-place. A second head can used to clear the original 1s, resulting in 

$$\begin{array}{c} \begin{array}{c} 1003\\ 1004\\ 1005\\ 1006\\ 1006\\ 1007\\ 1007\\ 1008\\$$

as desired. 

**Implementation of functions.** To illustrate a set of operations that could be implemented with a transformer, we consider approximating functions as sums of ReLUs; we use a result from Bai et al. (2023b), which we present below.

**Definition 1** (Definition 12 in Bai et al. (2023b)). A function  $q : \mathbb{R}^k \to \mathbb{R}$  is  $(\epsilon, R, M, C)$ -approximable by sum of ReLUs, if there exists an "(M,C)-sum of ReLUs" function 

1016  
1017 
$$f_{M,C}(\boldsymbol{x}) = \sum_{m=1}^{M} c_m ReLU(\boldsymbol{a}_m^{\top}[\boldsymbol{x};1]) \text{ with } \sum_{m=1}^{M} |c_m| \le C, \max_{m \in [M]} \|\boldsymbol{a}_m\|_1 \le 1, \ \boldsymbol{a}_m \in \mathbb{R}^{k+1}, \ c_m \in \mathbb{R}$$
1018

such that  $\sup_{\boldsymbol{x}\in [-R,R]^k} |g(\boldsymbol{x}) - f_{(M,C)}(\boldsymbol{x})| \leq \epsilon.$ 

**Definition 2** (Definition A.1 in Bai et al. (2023b)). We say a function  $g : \mathbb{R}^k \to \mathbb{R}$  is  $(R, C_l)$ -smooth if for  $s = \lceil (k-1)/2 \rceil + 2$ ,  $g \in C^{25}$  on  $[-R, R]^k$  and 

1022  
1023 
$$\sup_{\boldsymbol{x}\in[-R,R]^k} \left\|\nabla^i g(\boldsymbol{x})\right\|_{\infty} = \sup_{\boldsymbol{x}\in[-R,R]^k} \max_{j1,\dots,j_i\in[k]} \left|\partial_{x_{j1},\dots,x_{j_i}} g(\boldsymbol{x})\right| \le L_i$$

for all 
$$i = 0, 1, 2$$
 with  $\max_{0 \le i \le s} L_i R^i \le C_l$ .

 ${}^{5}C^{i}$  denotes that a function is *i* times differentiable with continuous *i*-th derivative.

**Proposition 1** (Proposition A.1 in Bai et al. (2023b)). For any  $\epsilon > 0$ ,  $R \ge 1$ ,  $C_l > 0$ , we have that: Any  $(R, C_l)$ -smooth function,  $g : \mathbb{R} \to \mathbb{R}$  is  $(\epsilon, R, M, C)$ -approximable by sum of ReLUs (Definition 1) with  $M \leq C(k)C_l^2 \log(1+C_l\epsilon)/\epsilon^2$ .

**Lemma 2.** For any function  $q : \mathbb{R}^k \to \mathbb{R}$  that is  $(R, C_l)$ -smooth, there exists a transformer with two layers, one head and width  $\mathcal{O}(\log(n) + d)$ , where d satisfies the requirements of Prop. 1, such that given as input 

$$\boldsymbol{X} = \begin{bmatrix} \boldsymbol{x}_1^{(1)} & \dots & \boldsymbol{y}^{(j-1)} & \boldsymbol{x}_1^{(j)} & \dots & \boldsymbol{x}_{n-2}^{(j)} & = & \boldsymbol{y}^{(j)} & \boldsymbol{x}_1^{(j+1)} & \dots \\ 0 & \dots & 1 & 0 & \dots & 0 & 0 & 1 & 0 & \dots \\ \boldsymbol{0} & \dots & \boldsymbol{0} & \boldsymbol{0} & \dots & \boldsymbol{0} & \boldsymbol{0} & \boldsymbol{0} & \dots \end{bmatrix} ,$$

it outputs

 $\boldsymbol{X} = \begin{bmatrix} \boldsymbol{x}_1^{(1)} & \dots & \boldsymbol{y}^{(j-1)} & \boldsymbol{x}_1^{(j)} & \dots & \boldsymbol{x}_{n-2}^{(j)} & = & \boldsymbol{y}^{(j)} & \boldsymbol{x}_1^{(j+1)} & \dots \\ * & \dots & * & * & \dots & * & * & * & * & \dots \\ \boldsymbol{0} & \dots & \tilde{g}(\boldsymbol{x}_i^{(j-1)}) - \boldsymbol{y}^{(j-1)} & * & \dots & * & * & \tilde{g}(\boldsymbol{x}_i^{(j)}) - \boldsymbol{y}^{(j)} & * & \dots \end{bmatrix}$ 

where 
$$|\tilde{g}(x) - g(x)| \leq \epsilon$$
 and for some i

*Proof.* We consider that the positional encodings are added in the input and we have

1044				-		-		•			
1045		$\begin{bmatrix} x_1^{(1)} \end{bmatrix}$		$oldsymbol{x}_1^{(j)}$	$oldsymbol{x}_2^{(j)}$		$x_{n-2}^{(j)}$	=	$oldsymbol{y}^{(j)}$	$oldsymbol{x}_1^{(j+1)}$	]
1046		0		0	0		0	0	1	0	
1047	X =	0		0	0		0	0	0	0	
1048		1	•••	1	1	• • •	1	1	1	1	
1049		$p_{n+1}$	• • •	$oldsymbol{p}_{jn+1}$	$p_{jn+2}$	• • •	$p_{jn+n-2}$	$oldsymbol{p}_{jn+n-1}$	$p_{jn+n}$	$p_{(j+1)n+1}$	• • •
1050		$p_{n+1-s}$	•••	$p_{jn+1-s}$	$p_{jn+2-s}$	• • •	$p_{jn+n-2-s}$	$p_{jn+n-1-s}$	$p_{jn+n-s}$	$p_{(j+1)n+1-s}$	•••_
1051	whom	wa fir co		aitional and	adinas m	<b>b</b>	m⊤m is long	an than m⊤m k	the same the	(6)	

where we fix some positional encodings  $p_k$  where  $p_k^{\top} p_k$  is larger than  $p_k^{\top} p_l$  by some threshold for  $k \neq l$ . The encodings used here are the binary representations of  $k \in \{-1, 1\}^{\log(mn)}$ . Further, we consider 1s in the positions with the results of the task to differentiate the context of the task and the result of the task. Define s = n - i, the distance between the result and the associated value in the context. 

In the first layer, we use the MLP's to create  $\tilde{q}$  according to Proposition 1 

 $\cdots$   $p_{jn+n-2-s}$   $p_{jn+n-1-s}$   $p_{jn+n-s}$  $p_{(i+1)n+1-s}$ 

The next operation is a shift of the sequence of  $\tilde{q}(\cdot)$ 's to the right by s. This will align the desired output  $\tilde{g}(\boldsymbol{x}_{i}^{(j)})$  with the observed output  $\boldsymbol{y}^{(j)}$ . Consider the following weight matrices 

$$\boldsymbol{W}_Q = \begin{bmatrix} \dots & 0 & \boldsymbol{\mathsf{I}} \end{bmatrix} \tag{8}$$

$$\boldsymbol{W}_{K} = \begin{bmatrix} \dots & 0 & C\mathbb{I} & \boldsymbol{0} \end{bmatrix}$$
(9)

$$\boldsymbol{W}_{V} = \begin{bmatrix} \boldsymbol{0} & \boldsymbol{0} & \boldsymbol{0} & \dots & \boldsymbol{0} \\ \boldsymbol{0} & \boldsymbol{0} & \boldsymbol{0} & \dots & \boldsymbol{0} \\ \boldsymbol{0} & \boldsymbol{0} & \boldsymbol{\mathbb{I}} & \dots & \boldsymbol{0} \\ \vdots & \vdots & \vdots & & \vdots \\ \boldsymbol{0} & \boldsymbol{0} & \boldsymbol{0} & \dots & \boldsymbol{0} \end{bmatrix}$$
(10)

for some large constant C to decrease error from the softmax attending to the incorrect tokens. This produces (within a small error induced by using a softmax)

 $(\boldsymbol{X}^{\top}\boldsymbol{W}_{K}^{\top}\boldsymbol{W}_{Q}\boldsymbol{X})_{i,j} = \boldsymbol{p}_{n+i}^{\top}\boldsymbol{p}_{n-s+j}$ (11) $\sigma_S(\boldsymbol{X}^{\top}\boldsymbol{W}_K^{\top}\boldsymbol{W}_Q\boldsymbol{X})_{i,j} = \mathbb{1}_{\{n+i=n-s+j\}} = \mathbb{1}_{\{i=j-s\}}$ (12) $\boldsymbol{W}_{V}\boldsymbol{X} = \begin{bmatrix} \cdots & \mathbf{0} & \mathbf{0} & \cdots & \mathbf{0} & \mathbf{0} & \mathbf{0} \\ \cdots & 0 & 0 & \cdots & 0 & 0 & 0 \\ \cdots & \tilde{g}(\boldsymbol{x}_{1}^{(j)}) & \tilde{g}(\boldsymbol{x}_{2}^{(j)}) & \cdots & \tilde{g}(\boldsymbol{x}_{n-2}^{(j)}) & \tilde{g}(\mathbf{z}) & \tilde{g}(\mathbf{z}_{n-2}^{(j)}) \\ \vdots & \vdots & \vdots & \vdots & \vdots \\ \cdots & \mathbf{0} & \mathbf{0} & \mathbf{0} & \mathbf{0} & \mathbf{0} & \mathbf{0} \end{bmatrix}$ . . . (13) $\boldsymbol{W}_{V}\boldsymbol{X}\sigma_{S}(\boldsymbol{X}^{\top}\boldsymbol{W}_{K}^{\top}\boldsymbol{W}_{Q}\boldsymbol{X}) = \begin{bmatrix} \cdots & \mathbf{0} & \mathbf{0} & \cdots & \mathbf{0} & \mathbf{0} & \mathbf{0} & \cdots & \cdots \\ \cdots & \mathbf{0} & \mathbf{0} & \cdots & \mathbf{0} & \mathbf{0} & \mathbf{0} & \cdots & \cdots \\ \cdots & * & * & \cdots & * & * & \tilde{g}(\boldsymbol{x}_{i}^{(j)}) & \cdots & \cdots \\ \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \cdots \\ \cdots & \mathbf{0} & \mathbf{0} & \cdots & \mathbf{0} & \mathbf{0} & \mathbf{0} & \cdots & \mathbf{0} \end{bmatrix}$   $\boldsymbol{X} + \boldsymbol{W}_{V}\boldsymbol{X}\sigma_{S}(\boldsymbol{X}^{\top}\boldsymbol{W}_{K}^{\top}\boldsymbol{W}_{Q}\boldsymbol{X}) = \begin{bmatrix} \cdots & \boldsymbol{x}_{1}^{(j)} & \boldsymbol{x}_{2}^{(j)} & \cdots & \boldsymbol{x}_{n-1}^{(j)} & \boldsymbol{z}_{n-1} & \boldsymbol{z} \\ \cdots & \mathbf{0} & \mathbf{0} & \cdots & \mathbf{0} & \mathbf{0} & \mathbf{0} \\ \cdots & * & * & \cdots & * & * & * \\ \cdots & \mathbf{0} & \mathbf{0} & \cdots & \mathbf{0} & \mathbf{0} & \mathbf{0} \\ \cdots & * & * & \cdots & * & * & * \\ \cdots & \mathbf{0} & \mathbf{0} & \cdots & \mathbf{0} & \mathbf{0} \\ \cdots & 1 & 1 & 1 & \cdots & 1 & 1 \\ \cdots & p_{jn+1} & p_{jn+2} & \cdots & p_{jn+n-2} & p_{jn+n-1} \\ \cdots & p_{jn+1-s} & p_{jn+2-s} & \cdots & p_{jn+n-2-s} & p_{jn+n-1-s} \end{bmatrix}$ (14). . . Ů . . . . . .  $p_{in+n}$  $p_{jn+2-s}$  ...  $p_{jn+n-2-s}$   $p_{jn+n-1-s}$   $p_{jn+n-s}$ (15) (16)Each matrix above only shows the slice that contains the j-th in-context example. This is repeated for each of the other in-context examples. As a final step with an MLP, subtract row 1 from row 3 to achieve the following output:  $\begin{bmatrix} \dots & \mathbf{x}_{1}^{(j)} & \mathbf{x}_{2}^{(j)} & \dots & \mathbf{x}_{n-1}^{(j)} & = & \mathbf{y}^{(j)} & \dots \\ \dots & 0 & 0 & \dots & 0 & 0 & 1 & \dots \\ \dots & * & * & \dots & * & * & \tilde{g}(\mathbf{x}_{i}^{(j)}) - \mathbf{y}^{(j)} & \dots \\ \dots & \mathbf{0} & \mathbf{0} & \dots & \mathbf{0} & \mathbf{0} & \mathbf{0} & \dots \\ \dots & 1 & 1 & \dots & 1 & 1 & 1 & \dots \\ \dots & \mathbf{p}_{jn+1} & \mathbf{p}_{jn+2} & \dots & \mathbf{p}_{jn+n-2} & \mathbf{p}_{jn+n-1} & \mathbf{p}_{jn+n} & \dots \\ \dots & \mathbf{p}_{jn+1} & \mathbf{p}_{jn+2} & \dots & \mathbf{p}_{jn+n-2} & \mathbf{p}_{jn+n-1} & \mathbf{p}_{jn+n} & \dots \end{bmatrix}$ (17) $p_{in+1-s}$  $p_{jn+2-s}$  ...  $p_{jn+n-2-s}$   $p_{jn+n-1-s}$  $p_{jn+n-s}$ **Copy Tasks** As has been experimentally investigated, the situation where a specific position within the context is copied as the label can be easily implemented by setting g(x) = x. The dependence on the subscript *i* within the construction is what allows the position copied to vary. D.2.1 IDENTIFYING IF TASK'S OUTPUT MATCHES THE IN-CONTEXT EXAMPLE **Lemma 3.** A three layer transformer with ReLU MLPs and embedding dimension  $\mathcal{O}(d + \log(mn))$ can calculate the proportion of in context examples that come from a specific task, where m is the number of in-context examples, each of length n and dimension d. 

*Proof.* We now have a matrix of the following form.

	]	$oldsymbol{y}^{(j)}$ .	=		$oldsymbol{x}_2^{(j)}$	$x_1^{(j)}$	<b>Г</b>
		1	0		Õ	Ō	
		$f(\boldsymbol{x}_{\cdot}^{(j)}) - \boldsymbol{y}^{(j)}$	*		*	*	
(18)		0	0		0	0	
		1 .	1		1	1	
		$p_{jn+n}$ .	$oldsymbol{p}_{jn+n-1}$		$p_{jn+2}$	$p_{jn+1}$	
	]	$p_{jn+n-s}$ .	$p_{jn+n-1-s}$	s	$p_{jn+2-2}$	$p_{jn+1-s}$	L
					<i></i>		
naxs and	rom sof	mall error coming f	, with some s	$^{j)} \approx 0$	$oldsymbol{x}_{\cdot}^{(j)}) - oldsymbol{y}_{\cdot}^{(j)}$	ct, than $f($	f the task is correc
MLP. For	sing an	of $f(\pmb{x}_{\cdot}^{(j)}) - \pmb{y}^{(j)}$ u	the L1-norm	e find	. First, we	ation erro	unction approxim
				use	$\boldsymbol{z}$ , we can	r arbitrary	calculating $\ m{z}\ _1$ fo
				d	_		
(19)		$\mathrm{U}(-oldsymbol{z}_i)$	$LU(\boldsymbol{z}_i) - ReL$	→ ]ReI	$\ z\ _{1} = $		
				=1	i		
			us wa hava	DTh	1 lover M	in a single	which can be done
			us, we have	LF, 111		in a single	
	]	$oldsymbol{u}^{(j)}$	=		$oldsymbol{x}_2^{(j)}$	$oldsymbol{x}_1^{(j)}$	Γ
		1	0		Õ	$\dot{0}$	
		$f(\boldsymbol{x}_{i}^{(j)}) - \boldsymbol{y}^{(j)}$	*		*	*	
(20)		$\ f(\boldsymbol{x}^{(j)}) - \boldsymbol{y}^{(j)}\ _{1}$	*		*	*	
(20)		$\  f(\boldsymbol{x}, \boldsymbol{y}) - \boldsymbol{y} \ _{1}$	0 0	•••	0	0	
		1	1		1	1	
		$p_{in+n}$	$p_{in+n-1}$		$p_{in+2}$	$oldsymbol{p}_{in+1}$	
	]	$p_{jn+n-s}$	$p_{jn+n-1-s}$		$p_{jn+2-s}$	$p_{jn+1-s}$	<b>.</b>
						4 1 . 1	
d be zero	ws wou	er task, the "extra" ro	on than anothe	nensio	ifferent di	task has c	Notice that if some
d be zero	ows wou	er task, the "extra" ro	on than anothe	nensio	ifferent di	task has d he result.	Notice that if some nd will not affect t
d be zero following	bws wou $\delta$ in the	er task, the "extra" ro	on than anothe	nensio row to	Ifferent different differ	task has d he result.	Notice that if some and will not affect t For clarity, we set a
d be zero following at a later	bws wou $\xi \hat{\delta}$ in the sh value	er task, the "extra" ro Il cause the following ion handles these tra	on than anothe 1s. These wi the construct	row to tted as	If the $\ \cdot\ _1$ on the $\ \cdot\ _1$	task has d he result. ll * values operation	Notice that if some and will not affect the For clarity, we set a set these to 0. This
d be zero following at a later	ws wou $_{3}\hat{\delta}$ in the sh value	er task, the "extra" ro Il cause the following ion handles these tra	on than anothe 1s. These wi the construct	nensio row to tted as	In the $\ \cdot\ _1$ can be omi	task has of he result. ll * values operation	Notice that if some and will not affect t For clarity, we set a set these to 0. This ayer.
d be zero following at a later	bws wou $\hat{\delta}$ in the sh value let x real	er task, the "extra" ro Il cause the following ion handles these tra	on than anothe 1s. These wi the construct	row to tted as	Ifferent due in the $\ \cdot\ _1$ can be omi	task has of he result. Il * values operation value of th	lotice that if some nd will not affect to or clarity, we set a et these to 0. This yer.
d be zero following at a later resent the	bws wou $\hat{\delta}$ in the sh value let x rep to 1	er task, the "extra" roll cause the following ion handles these training the $y$ vectors and thus set the $*$ values	on than anothe 1s. These wi the construct nd row markin following Re	row to tted as e seco	ifferent during the second du	task has c he result.    * values operation value of th ith $   f(x^{(j)})$	Notice that if some nd will not affect to for clarity, we set a et these to 0. This ayer. Let b represent the alues in the row w
d be zero following at a later resent the	bws wou $\hat{\delta}$ in the sh value let x rep to 1.	er task, the "extra" roll cause the following ion handles these training the $y$ vectors and LUs set the $*$ values	on than anothe o 1s. These wi the construct nd row markin following Re	row to tted as e seco 1. The	In the $\ \cdot\ _1$ can be omi e flag in th $(x^{(j)}) - y^{(j)}\ _1^2$	task has compared to the result. Il * values operation value of the ith $  f(\mathbf{x}^{(j)}) $	Notice that if some nd will not affect to for clarity, we set a et these to 0. This ayer. Let b represent the alues in the row w
d be zero following at a later resent the	bws wou $\hat{\delta}$ in the sh value let x rep to 1.	er task, the "extra" roll cause the following ion handles these training the $y$ vectors and LUs set the * values	on than another o 1s. These wi the construct nd row markin following Re	row to tted as e seco 1. The	ifferent different different different different different difference of the second s	task has constrained by the result. Il * values operation value of the second secon	Notice that if some nd will not affect to for clarity, we set a tet these to 0. This ayer. Let $b$ represent the alues in the row w
d be zero following at a later resent the (21)	we would be a constraint of $\hat{\delta}$ in the sh value let $x$ replaced to 1.	er task, the "extra" ro Il cause the following ion handles these tra ng the $y$ vectors and LUs set the * values U(Cb - C + 1)	on than another o 1s. These with the construct and row marking following Ref - Cb) - ReL	nension row to tted as e secon 1. The LU(x + 1)	ifferent different different different different different different difference of the second secon	task has consistent of the result. Il * values operation value of the inth $  f(x)  \leq x \leq x$	fotice that if some and will not affect to or clarity, we set a tet these to 0. This ayer. b represent the alues in the row w
d be zero following at a later esent the (21) s reduces	we would be a constraint of the second seco	For task, the "extra" roll cause the following ion handles these training the $y$ vectors and LUs set the * values U(Cb - C + 1) -x = 1 and when	on than another o 1s. These with the construct and row marking following Ref - Cb) - ReL	nension row to tted as e seco t. The LU(x + t)	In the $\ \cdot\ _1$ can be omised for the flag in the flag	task has c he result. Il * values operation value of th ith $  f(\mathbf{x}^{(x)} \times \mathbf{x} \leftarrow \mathbf{x}) $ $x \leftarrow \mathbf{x}$ tant C W	Notice that if some nd will not affect to for clarity, we set a et these to 0. This hyer. et <i>b</i> represent the alues in the row w
d be zero following at a later resent the (21) s reduces	we would be a constraint of the second seco	er task, the "extra" roll cause the following ion handles these training the $y$ vectors and LUs set the * values U(Cb - C + 1) -x = 1, and when a	on than another o 1s. These with the construct and row marking following Ref - Cb) - ReL duces to $x + 1$	nension row to tted as e second. The LU(x + this real	inferent different differ	task has c he result.    * values operation value of th ith $   f(x^{(j)} \times x \leftarrow x)$ $x \leftarrow x$ tant C. W as desired.	Notice that if some and will not affect to For clarity, we set a set these to 0. This ayer. Let <i>b</i> represent the values in the row w
d be zero following at a later resent the (21) s reduces	we would be a constraint of $\delta$ in the sh value let $x$ replaced to 1. b = 1, the second	er task, the "extra" roll cause the following ion handles these training the $y$ vectors and LUs set the * values U(Cb - C + 1) -x = 1, and when a	on than another o 1s. These with the construct and row marking following Ref - Cb) - ReL duces to $x + 1$	nension row to tted as e secont. The LU(x + this real	ifferent different differ	task has c he result. Il * values operation value of th ith $  f(x^{(x)} + x - x) $ $x \leftarrow x$ tant C. Was desired.	Notice that if some and will not affect to For clarity, we set a set these to 0. This layer. Let <i>b</i> represent the values in the row we for some large cons to $x + 1 - 1 = x$ , a
d be zero following at a later resent the (21) s reduces	we would be a constrained by a constraint of the set o	er task, the "extra" ro Il cause the following ion handles these tra ing the $y$ vectors and LUs set the * values U(Cb - C + 1) -x = 1, and when a	on than another o 1s. These with the construct and row marking following Ref - Cb) - ReL duces to $x + 1$	nension row to tted as e seco t. The LU(x + this recommended)	in the $\ \cdot\ _1$ can be omised flag in the	task has c he result. Il * values operation value of th ith $  f(x^{(j)} \times c_x) $ $x \leftarrow x$ tant C. W is desired. $x^{(j)}$	Notice that if some and will not affect to For clarity, we set a set these to 0. This ayer. Let <i>b</i> represent the values in the row we for some large cons to $x + 1 - 1 = x$ , a
d be zero following at a later resent the (21) s reduces	we would be a constraint of $\hat{\delta}$ in the sh value let $x$ repload to 1. b = 1, the second secon	er task, the "extra" roll cause the following ion handles these training the $y$ vectors and LUs set the * values U(Cb - C + 1) $-x = 1$ , and when $y_1^{(j)}$	on than another on the construct the construct and row marking following Ref - Cb) - ReL luces to $x + 1$	nension row to tted as e secont. The LU( $x$ - this red	inferent different different different different different different difference of the second secon	task has c he result. Il * values operation value of th ith $  f(x^{(j)}) $ $x \leftarrow a$ tant C. We as desired. $x_1^{(j)}$	Notice that if some and will not affect to For clarity, we set a set these to 0. This ayer. Let <i>b</i> represent the values in the row w For some large cons to $x + 1 - 1 = x$ , a
d be zero following at a later esent the (21) s reduces	we would be a constraint of the second seco	er task, the "extra" ro Il cause the following ion handles these tra ing the $y$ vectors and LUs set the * values U(Cb - C + 1) -x = 1, and when $yf(x_{i}^{(j)}) (i)$	on than another on the construct the construct and row markin following Ref - Cb) - ReL duces to $x + 1$ = 0	nension row to tted as e seco $_1$ . The LU(x + this red	inferent different different different different different different difference of the second secon	task has c he result. Il * values operation value of th ith $  f(x^{(j)} + c$ $x \leftarrow a$ tant C. We tast desired. $x_1^{(j)}$ 0	Notice that if some and will not affect to For clarity, we set a set these to 0. This layer. Let <i>b</i> represent the values in the row w for some large cons to $x + 1 - 1 = x$ , a $\begin{bmatrix} \dots \\ \dots \end{bmatrix}$
d be zero following at a later resent the (21) s reduces	we would be a constraint of the second seco	er task, the "extra" ro Il cause the following ion handles these tra ing the $y$ vectors and LUs set the * values U(Cb - C + 1) -x = 1, and when a $y^{(j)}_{1}$ $f(x^{(j)}_{2}) - y^{(j)}_{2}$	on than another on the construct the construct and row markin following Ref - Cb) - ReL duces to $x + 1$ = 0	nension row to tted as e secont. The LU(x + this red) 	inferent different differ	task has c he result. Il * values operation value of th ith $  f(x^{(j)} + x \leftarrow x) $ tant C. W tant C. W tas desired. $x_1^{(j)}$ 0 *	Notice that if some and will not affect to For clarity, we set a set these to 0. This layer. Let <i>b</i> represent the values in the row w for some large cons to $x + 1 - 1 = x$ , a $\begin{bmatrix} \dots \\ \dots \\ \dots \\ \dots \end{bmatrix}$
d be zero following at a later esent the (21) s reduces (22)	by swore $\hat{\delta}$ in the sh value let $x$ reproduces to 1. b = 1, the second secon	er task, the "extra" roll cause the following ion handles these training the $\boldsymbol{y}$ vectors and LUs set the * values U(Cb - C + 1) $-x = 1$ , and when $y_{1}^{(j)}$ $f(\boldsymbol{x}_{\cdot}^{(j)}) - \boldsymbol{y}^{(j)}$ $\ f(\boldsymbol{x}_{\cdot}^{(j)}) - \boldsymbol{y}^{(j)}\ _{1}$	on than another on the construct of the construct and row marking following Ref - Cb) - ReL huces to $x + 1$ = 0 * 1	row to row to tted as e seco t. The LU(x this red	inferent different differ	task has c he result. Il * values operation value of th ith $  f(x^{(j)} + x \leftarrow x) $ tant C. Was desired. $x_1^{(j)} = 0$ * 1	Notice that if some nd will not affect to For clarity, we set a et these to 0. This ayer. Let <i>b</i> represent the alues in the row w or some large cons to $x + 1 - 1 = x$ , a $\left[ \begin{array}{c} \dots \\ \dots \\ \dots \\ \dots \\ \dots \end{array} \right]$
d be zero following at a later resent the (21) s reduces (22)	by swor $\hat{\delta}$ in the sh value let $x$ rep to 1. $b = 1$ , the short $\hat{\delta}$ is the short $\hat{\delta}$ in the short $\hat{\delta}$ in the short $\hat{\delta}$ in the short $\hat{\delta}$ in the short $\hat{\delta}$ in the short $\hat{\delta}$ in the short $\hat{\delta}$ in the short $\hat{\delta}$ in the short $\hat{\delta}$ is the short $\hat{\delta}$ in the short $\hat{\delta}$ is the short $\hat{\delta}$ in the short $\hat{\delta}$ in the short $\hat{\delta}$ in the short $\hat{\delta}$ in the short $\hat{\delta}$ in the short $\hat{\delta}$ in the short $\hat{\delta}$ in the short $\hat{\delta}$ in the short $\hat{\delta}$ in the short $\hat{\delta}$ in the short $\hat{\delta}$ in the short $\hat{\delta}$ in the short $\hat{\delta}$ in the short $\hat{\delta}$ in the short $\hat{\delta}$ in the short $\hat{\delta}$ in the short $\hat{\delta}$ is the short $\hat{\delta}$ in the short $\hat{\delta}$ is the	er task, the "extra" roll cause the following ion handles these training the $y$ vectors and LUs set the * values U(Cb - C + 1) $-x = 1$ , and when $y_{1}^{(j)}$ $f(x_{.}^{(j)}) - y_{.}^{(j)}$ $\ f(x_{.}^{(j)}) - y_{.}^{(j)}\ _{1}$ 0	on than another on the construct the construct of row marking following Ref -Cb) - ReLluces to x + 1 = 0 * 1 0	mension row to tted as e seco $_1$ . The LU(x - this red	inferent due in the $\ \cdot\ _1$ can be omining the flag in the $ \cdot  _1 - \mathbf{R} = \mathbf{R}$ $\mathbf{r}_2^{(j)} = \mathbf{r}_2^{(j)}$ $\mathbf{r}_2^{(j)} = \mathbf{r}_2^{(j)}$ $\mathbf{r}$	task has c he result. Il * values operation value of th ith $  f(x^{(j)} + x \leftarrow x) $ tant C. We tast desired. $x_1^{(j)} = 0$ * 1 0	Notice that if some nd will not affect to For clarity, we set a et these to 0. This ayer. Let <i>b</i> represent the alues in the row w
d be zero following at a later resent the (21) s reduces (22)	by swore we have $\delta$ in the sh value let $x$ represented by the shift $b = 1$ , the shift $b = 1$ , the shift $b = 1$ , the shift $b = 1$ , the shift $b = 1$ , the shift $b = 1$ , the shift $b = 1$ , the shift $b = 1$ , the shift $b = 1$ and $b = 1$ .	er task, the "extra" roll cause the following ion handles these training the $y$ vectors and LUs set the * values U(Cb - C + 1) $-x = 1$ , and when $y_{j}^{(j)}$ $f(\boldsymbol{x}_{\cdot}^{(j)}) - \boldsymbol{y}_{j}^{(j)}$ $\ f(\boldsymbol{x}_{\cdot}^{(j)}) - \boldsymbol{y}_{\cdot}^{(j)}\ _{1}$ 0	on than another on than another on the construct the construct and row marking following Ref -Cb) - ReLhuces to $x + 1=0*101$	mension row to tted as e seco t. The LU( $x$ this red this red this red	ifferent different differ	task has c he result. Il * values operation value of th ith $  f(x^{(j)} + x \leftarrow x) $ tant C. W is desired. $x_1^{(j)} = 0$ $x_1^{(j)} = 0$ 1 0	Notice that if some nd will not affect the For clarity, we set a tet these to 0. This ayer. Let b represent the alues in the row we for some large cons to $x + 1 - 1 = x$ , a x + 1 - 1 = x, and $x + 1 - 1 = x$ , and $x + 1 - 1 = x$ .
d be zero following at a later resent the (21) s reduces (22)	by swore we have $\delta$ in the sh value let $x$ represented by the shift $b = 1$ , the shift $b = 1$ and the shift $b = 1$ .	er task, the "extra" roll ll cause the following ion handles these tra ing the $y$ vectors and LUs set the * values U(Cb - C + 1) -x = 1, and when $af(\boldsymbol{x}_{\cdot}^{(j)}) - \boldsymbol{y}^{(j)}\ f(\boldsymbol{x}_{\cdot}^{(j)}) - \boldsymbol{y}^{(j)}\ _{1}01p_{jn+n}p_{jn+n}$	on than another on than another on the construct the construct and row marking following Refined - Cb) - ReLind huces to $x + 1= 0*101p_{jn+n-1}$	mension row to tted as e seco $_1$ . The LU(x - this red)   	inferent different differ	task has c he result. Il * values operation value of th ith $  f(x^{(j)} + x \leftarrow a) $ tant C. With tant C. With tant C. With $x_1^{(j)} = 0$ * 1 $p_{jn+1}$	otice that if some and will not affect to or clarity, we set a set these to 0. This yer. et <i>b</i> represent the ilues in the row w

Now define a thresholding function  $\delta(z)$  that satisfies  $\delta(0) = 1$  and  $\delta(z) = 0$  for z >> 0. One such function used here is

$$\hat{\delta}_C(z) = \text{ReLU}(1 - Cz) \tag{23}$$

for some constant C, where larger C captures a narrower neighborhood of 0.

However, a slight change needs to be added to  $\hat{\delta}_C$ . In the same row as  $||f(x^{(j)}) - y^{(j)}||$  are many values that need to be discarded. Let b be the bit for the current column marking if the column contains an x or a y. We use instead

 $\hat{\delta}_C(b,z) = \text{ReLU}(b-Cz)$ (24)This will be zero whenever b = 0 and z > 0. We then have as output  $\begin{bmatrix} \dots & \boldsymbol{x}_{1}^{(j)} & \boldsymbol{x}_{2}^{(j)} & \dots & = & \boldsymbol{y}^{(j)} & \dots \\ \dots & 0 & 0 & \dots & 0 & 1 & \dots \\ \dots & * & * & \dots & * & f(\boldsymbol{x}_{\cdot}^{(j)}) - \boldsymbol{y}^{(j)} & \dots \\ \dots & 1 & 1 & \dots & 1 & \|f(\boldsymbol{x}_{\cdot}^{(j)}) - \boldsymbol{y}^{(j)}\|_{1} & \dots \\ \dots & 0 & 0 & \dots & 0 & \hat{\delta}_{C}(\|f(\boldsymbol{x}_{\cdot}^{(j)}) - \boldsymbol{y}^{(j)}\|_{1}) & \dots \\ \dots & 0 & \mathbf{0} & \dots & \mathbf{0} & \mathbf{0} & \dots \\ \dots & 1 & 1 & \dots & 1 & 1 & \dots \\ \dots & 1 & 1 & \dots & 1 & 1 & \dots \end{bmatrix}$ (25) $p_{jn+1}$   $p_{jn+2}$   $\ldots$   $p_{jn+n-1}$  $p_{jn+n}$  $p_{jn+1-s}$   $p_{jn+2-s}$  ...  $p_{jn+n-1-s}$  $p_{jn+n-s}$ 

1213 Importantly,  $\hat{\delta}_C(\|f(\boldsymbol{x}_{\cdot}^{(j)}) - \boldsymbol{y}^{(j)}\|_1) = 1$  when  $f(\cdot)$  is the correct task and  $\hat{\delta}_C(\|f(\boldsymbol{x}_{\cdot}^{(j)}) - \boldsymbol{y}^{(j)}\|_1) = 0$ 1214 when  $f(\cdot)$  disagrees by more than  $\frac{1}{C}$  in L1-norm.

Lastly, for the next step in the construction, we need to average these soft indicators  $\hat{\delta}$  to see how 1217 common f is within the context. This is done with an attention layer. Let  $W_Q$  select the row with all 1218 1218 

$$\boldsymbol{X}^{\top} \boldsymbol{W}_{K}^{\top} \boldsymbol{W}_{Q} \boldsymbol{X} = \begin{bmatrix} \vdots & \vdots & \vdots & \vdots & \vdots \\ 0 & 0 & \dots & 0 & 0 \\ 0 & 0 & \dots & 0 & 0 \\ \vdots & \vdots & \ddots & \vdots & \vdots \\ 0 & 0 & \dots & 0 & 0 \\ C & C & \dots & C & C \\ \vdots & \vdots & & \vdots & \vdots \end{bmatrix}$$
(26)  
$$\sigma_{S}(\boldsymbol{X}^{\top} \boldsymbol{W}_{K}^{\top} \boldsymbol{W}_{Q} \boldsymbol{X}) \approx \begin{bmatrix} \vdots & \vdots & \vdots & \vdots & \vdots \\ 0 & 0 & \dots & 0 & 0 \\ 0 & 0 & \dots & 0 & 0 \\ \vdots & \vdots & & \vdots & \vdots & \vdots \\ 0 & 0 & \dots & 0 & 0 \\ 1/m & 1/m & \dots & 1/m & 1/m \\ \vdots & \vdots & & \vdots & \vdots & \vdots \end{bmatrix}$$
(27)

where a 1/m will appear in every row corresponding to a result  $\boldsymbol{y}$ . Let the value matrix select the row containing  $\hat{\delta}(\|f(\boldsymbol{x}^{(j)} - \boldsymbol{y}^{(j)})\|_1)$ . Denote  $p = \frac{1}{m} \sum_{j=1}^m \hat{\delta}(\|f(\boldsymbol{x}^{(j)} - \boldsymbol{y}^{(j)})\|_1)$ . Without causal masking, we would have as output

1242													
1243		Г	$\sigma^{(j)}$	<i>m</i> <sup>(</sup>	(j)		_			a(i)		٦	
1244			$x_1 \\ 0$	$x_2$	2	• • •	=			<b>y</b> <sup>(3)</sup>			
1245			0	C	)	•••	0		C (	(j)	(i)		
1246			*	k	k	• • •	*		$f(\boldsymbol{x})$	(i) -	$-y^{(j)}$		
1247			1	1	L	•••	1		$\ f(\boldsymbol{x}$	( <i>J</i> ) –	$oldsymbol{y}^{(j)} \Vert_1$		
1248			0	(	)		0	$\hat{\delta}$	$C(\ f(z)\ )$	$x_{\cdot}^{(j)})$ .	$- y^{(j)} \ _1)$		(28)
12/10			p	Į	)		p			p			
1245			0	(	)	•••	0			0			
1250			1	1	L		1			1			
1251			$oldsymbol{p}_{jn+1}$	$p_{jn}$	i+2	•••	$p_{jn+n}$	$^{-1}$		$p_{jn+}$	n		
1052		[···	$p_{jn+1-s}$	$p_{jn+}$	-2-s	$\cdots p$	$p_{jn+n-}$	1-s	1	$o_{jn+n}$	-s	· · · ]	
1255	However	with c	aucal macl	cing v	ve can	only g	uarante	e that r	o will a	nnear	in the colu	imns coi	ntaining
1254	the most r	ecent of	example be	eing a	ueried	. Thank	fully, f	his is al	l that i	s need	led.		
1200	uie most i	ceent	enumpie or		aerrea	. I mum	iiuiij, t	1110 10 <b>u</b>	ii tilut I	5 need			
1230		sk Ev	ECUTION										
1257	D.5 1A	SK LA	ECUTION										
1258	Lemma 4	. A tw	o laver trai	nsform	ier, wi	th embe	dding d	limensi	on O(a	$l + \log l$	g(mn)) ca	n perfori	n a task
1259	and weigh	nt its o	utput by the	e prop	ortion	of exa	mples o	f that to	ask see	n with	in the con	text.	
1260	0		1 5			5	1						
1261	Now that	the pro	oportions o	f each	task l	nave be	en iden	tified in	n the co	ontext	, the task i	tself nee	ds to be
1262	executed f	for the	new exam	ple be	ing qu	eried.	Го simp	olify not	tation,	let the	e input to the	nis step ł	be
1263													
1264					Γ	$x_1^{(m)}$	$x_2^{(m)}$		$x_{-2}^{(m)}$	=]			
1265						*	*		*	*			
1266				$\mathbf{V}$ –		p	p		p	p			(20)
1267				$\Lambda$ –		0	0		0	0			(29)
1268						*	*		*	*			
1269					[	1	1		1	1			
1270	<b>D</b> .11	41			41.1.1		. 1. 1	1 24	1.1.1.4	1 1.00		• • • • • • •	
1271	Following	the sa	ime proces	s as oi	ltimed	above,	, althou	gn with	i slight		erent posit	ional end	codings,
1272	calculate j	f(x)	) and plac	e that	result	in the r	inal col	umn be	ing dec		I nese nee	to be a	added at
1273	the beginn	ing of	the consu	uction	i, but a	ue omy	muou	uceu ne		Janty	•		
1274				-						\ <b>-</b>			
1275					$x_1^{(m)}$	$x_2^{(m)}$	,	$x_{n-2}^{(m)}$	$x_n^{(r)}$	$\binom{n}{-1}$			
1276					*	*		*	*				
1277					p	p		p	p				
1278					*	*		*	$f(oldsymbol{x}^{0})$	$^{(m)})$			$\langle 2 0 \rangle$
1279					0	0	• • •	0	C	)			(30)
1280					*	*	• • •	*	*	•			
1281				• • •	1	1		1	1	.			
1282					0	0	• • •	0					
1283				L···	0	0	• • •	1	U	L ۱			
1284	We will tr	ansfor	m the row o	ontair	ning al	1 n to h	e ahle t	o annro	vimate	lv mu	ltinly <i>n</i> by	$f(\mathbf{r}^{(m)})$	Using
1285	the second	to las	t row perfe	$\operatorname{rm} n$	$\rightarrow 1$ -	-n Us	ing the	last two	rows	clear	out the res	f of that	row and
1286	fill it with	-C for	or some lar	ge co	nstant	C. We	then ha	ive	, 10, 10,	ereur	out the res	t of that i	iow und
1287				8									
1288				г	(m)	(m)	)	(m)	(1	<i>n</i> ) 7			
1289				· · ·	$x_1^{(1,5)}$	$x_2$	•••	$x_{n-2}$	$oldsymbol{x}_n^{\scriptscriptstyle (\cdot)}$	_1			
1290				· · ·	*	*	•••	*	*	· _			
1291				· · ·	-C	-0		p	1 - r ( (	$\begin{bmatrix} p \\ m \end{bmatrix}$			
1292				· · ·	*	*	•••	*	$J(\boldsymbol{x})$				(31)
1293				· · ·	U	U		U	U.	.			(31)
1294					* 1	* 1		* 1	*	•			
1295					0	1	• • •	0	1	•			
1233					0	0	• • •	1	1				
				Γ	0	0	• • •	-	U	· _			

Further, use the second-to-last row to clear out all \* in the rows below.

$$\begin{bmatrix} \dots & \boldsymbol{x}_{1}^{(m)} & \boldsymbol{x}_{2}^{(m)} & \dots & \boldsymbol{x}_{n-2}^{(m)} & \boldsymbol{x}_{n-1}^{(m)} \\ \dots & * & * & \dots & * & * \\ \dots & -C & -C & \dots & p & 1-p \\ \dots & \mathbf{0} & \mathbf{0} & \dots & \mathbf{0} & f(\boldsymbol{x}^{(m)}) \\ \dots & \mathbf{0} & \mathbf{0} & \dots & \mathbf{0} & \mathbf{0} \\ \dots & * & * & \dots & * & * \\ \dots & 1 & 1 & \dots & 1 & 1 \\ \dots & 0 & 0 & \dots & 0 & 1 \\ \dots & 0 & 0 & \dots & 1 & 0 \end{bmatrix}$$
(32)

These previous operations can all be done in a single MLP.

1310 Lastly, use an attention layer where  $W_K$  selects the row with the -Cs,  $W_Q$  selects the row with all 1311 1s, and  $W_V$  selects the  $f(\mathbf{x}^{(m)})$ . For the last token  $\mathbf{x}_L$ ,

$$\boldsymbol{X}^{\top} \boldsymbol{W}_{K}^{\top} \boldsymbol{W}_{Q} \boldsymbol{x}_{L} = \begin{bmatrix} \vdots \\ -C \\ -C \\ \vdots \\ p \\ 1-p \end{bmatrix} \begin{bmatrix} 1 \end{bmatrix} = \begin{bmatrix} \vdots \\ -C \\ -C \\ \vdots \\ p \\ 1-p \end{bmatrix}$$
(33)  
$$\sigma_{S}(\boldsymbol{X}^{\top} \boldsymbol{W}_{K}^{\top} \boldsymbol{W}_{Q} \boldsymbol{x}_{L}) \approx \begin{bmatrix} \vdots \\ -\infty \\ -\infty \\ \vdots \\ p \\ 1-p \end{bmatrix} = \begin{bmatrix} \vdots \\ 0 \\ 0 \\ \vdots \\ \frac{1}{1+e^{1-2p}} \\ 1-\frac{1}{1+e^{1-2p}} \end{bmatrix}$$
(34)

$$\boldsymbol{w}_{V}\boldsymbol{x}_{L}\sigma_{S}(\boldsymbol{X}^{\top}\boldsymbol{W}_{K}^{\top}\boldsymbol{W}_{Q}\boldsymbol{x}_{L}) = \begin{bmatrix} \vdots \\ \mathbf{0} \\ \frac{1}{1+e^{1-2p}}\mathbf{0} + (1 - \frac{1}{1+e^{1-2p}})f(\boldsymbol{x}^{(m)}) \\ \mathbf{0} \\ \vdots \end{bmatrix}$$
(35)  
$$\boldsymbol{x}_{L} + \boldsymbol{W}_{V}\boldsymbol{x}_{L}\sigma_{S}(\boldsymbol{X}^{\top}\boldsymbol{W}_{K}^{\top}\boldsymbol{W}_{Q}\boldsymbol{x}_{L}) = \begin{bmatrix} \boldsymbol{x}_{n-1}^{(m)} \\ * \\ 1 - p \\ \frac{1}{1+e^{1-2p}}f(\boldsymbol{x}^{(m)}) \\ \mathbf{0} \\ * \\ 1 \\ 0 \end{bmatrix}$$
(36)

1347 Importantly, we are left with  $\frac{1}{1+e^{1-2p}}f(\boldsymbol{x}^{(m)})$ . The factor  $\frac{1}{1+e^{1-2p}}$  is approximately p, especially 1348 around  $\frac{1}{2}$ . This multiplication can also be calculated more accurately with approximations using 1349 ReLUs or sigmoids, but for brevity and following experimental evidence of a sigmoid shape in task superpositions, these options are ommited.

# 1350 D.4 SUPERPOSED TASKS WITH PARALLEL HEADS

The above construction works for a single task, where the output is weighted by the proportions of
 the task within the context. To complete the construction of a transformer that does superposition of
 tasks, each of these models needs to be placed within the same overall transformer. This is described
 here.

Let there be a collection of tasks  $\{t_i\}_{i=1}^T$  which can be executed by transformers with model weights represented by subscripts  $(\cdot_i)$ . With the input to each transformer being  $X^{(i)}$ , the overall input matrix is given by vertically stacking these matrices.

 $egin{array}{lll} egin{array}{c} egin{arra$ 

Similarly, define each MLP's weights and biases as

$$\boldsymbol{W} = \operatorname{diag}(\boldsymbol{W}_1, \dots, \boldsymbol{W}_T) \quad \boldsymbol{b} = \begin{bmatrix} \boldsymbol{b}_1 \\ \vdots \\ \boldsymbol{b}_T \end{bmatrix}$$
(38)

(37)

This puts every MLP to be independent of each other. Lastly, we need to change the attention layers. This requires the use of one head per task. In each of the following,  $W^{(i)}$  is a weight matrix for head  $i, (W)_i$  is the weight matrix for task i in its individual transformer, and each matrix below is in the *i*-th block.

$$\boldsymbol{W}_{V}^{(i)} = \begin{bmatrix} \vdots \\ \mathbf{0} \\ (\boldsymbol{W}_{V})_{i} \\ \vdots \end{bmatrix}^{\top} \quad \boldsymbol{W}_{K}^{(i)} = \begin{bmatrix} \vdots \\ \mathbf{0} \\ (\boldsymbol{W}_{K})_{i} \\ \vdots \end{bmatrix}^{\top} \quad \boldsymbol{W}_{Q}^{(i)} = \begin{bmatrix} \vdots \\ \mathbf{0} \\ (\boldsymbol{W}_{Q})_{i} \\ \vdots \end{bmatrix}^{\top}$$
(39)

In all, this model executes multiple tasks in superposition by using parallel streams of heads that each performs a single task. Task identification can happen through the same mechanism as task execution by comparing the output of the task on each in context example with the true output.

For context related tasks, there needs to be positional encodings that allow for looking back a fixed number of tokens. For context agnostic tasks, a wide MLP can be used to approximate arbitrary non-linear transformations of the input. Each of these tasks only require a small number of layers, significantly smaller than those of modern LLMs. It may be possible that LLMs do certain tasks with different combinations of layers.

Also, if we take the feature p from each parallel stream, this creates the following task identifier.

Interpolating between the pure tasks, represented by unit vectors, different amounts of each task will appear in the superposition in roughly equal proportions to those found in v.

Lastly, we restate this construction formally.

**Theorem 1.** A seven layer transformer with embedding dimension  $O(d + \log(mn))$  with K heads per attention layer can perform k tasks on vectors of dimension d in superposition, with weighting based on m different in-context examples each of length n.

*Proof.* Using in succession each of Lemma 1, Lemma 2, Lemma 3, and Lemma 4, a transformer with1409the desired properties can execute k tasks in parallel. Lemma 1 identifies positions within the context1410that contain the labels y. Lemma 2 then uses function approximation to perform arbitrary tasks1411within the architecture, which are then used by 3 to find the proportions of each task and aggregate1412them into a single task identifier. Lastly, Lemma 4 uses this task identifier to create a weighted sum1413of outputs from the different tasks based on their in-context proportions.

Remark. Transformers of greater depth than seven layers can also represent this construction by setting the weights in all other layers for the non residual part to zero.



For the non-zero probabilities we observe on each task answer, are they indications of task superposition or by-products of prediction noise?

1510

1512 We set up an experiment to investigate this. For each setting, we select two tasks task1 and task2; 1513 then we consider two scenarios: (1) we provide prompts where all task examples come from task1 1514 and (2) we provide prompts where half of the task examples come from task1 and the other half of 1515 task examples come from task2; in both scenarios we measure the probabilities for task answers of 1516 task1 and task2 and see how these probabilities change between scenario (1) and (2). For each scenario, we test it on 100 prompts (each task has 10 in-context examples) and plot the median of 1517 task answers in Figure 7. As is shown in Figure 7 left side, where there is no task example from 1518 task2 in the prompt, the probabilities for task answers of task2 are near 0 for all models; on the 1519 right side, where there is an equal number of task examples that are from task1 and task2, the 1520 probabilities for task answers of task2 increase significantly. This indicates that when we provide 1521 prompts that mix task examples from different tasks, the prediction on each task answer is more than 1522 just pure prediction noise. 1523

1523 1524 1525

## E.2 MORE ANALYSIS ON OTHER CATEGORY IN MODEL'S OUTPUT DISTRIBUTION

In Figure 2, in settings 1, 2 and 4 (correspond to Figure 2a, 2b and 2d respectively), since there are also non-negligible probabilities on the other category (and it is a summation of all other probabilities), we further investigate the probabilities in other category. In particular, for each prompt, we use beam search (where we stop searching when we encounter "\n") to find answers that have top-(K + 1) probabilities and record the maximum probability of the answer that is not one of the task answers. In Figure 8, we plot the median of such probability along with medians of probabilities of each task answer.



Figure 8: For each subplot, we plot the medians of probabilities for each task answer and the medianof the maximum probability of the answer that is not one of the task answers.

In Figure 8a (setting 1), for GPT-3.5 and Llama-3, we can observe that the medians of the probabilities
 for most probable non-task-answer are significantly lower than that of each task answer. This indicates
 that the models are effectively performing task superposition across the four tasks, with any other

answer having a low probability. For Qwen-1.5, however, the median probabilities of task answers for add, add\_in\_fr, add\_in\_es and the most probable non-task-answer are relatively low. This may be attributed to the model's limited ability to perform the add\_in\_fr and add\_in\_es tasks.

1569 In Figure 8b (setting 2), we see that for all models, the median probability of the most probable 1570 non-task-answer is lower than that of each task answer. A possible explanation for the most probable 1571 non-task-answer still having a probability around 0.07 is the presence of related tasks: (1) CAPITAL, 1572 which returns the capital of a country in uppercase, (2) CONTINENT, which returns the continent of 1573 a country in uppercase, and (3) identity, which directly returns the country name. If these task 1574 answers are excluded when calculating the top-(K + 1) probabilities, and the maximum probability 1575 of the non-task-answer is recalculated, the median probability drops to less than 0.02 for Qwen-1.5 1576 and Llama-3, and less than 0.001 for GPT-3.5. This suggests that while the models may not always strictly perform task superposition on the provided tasks, they may perform compositions of the 1577 tasks. We think it is an interesting direction for the future work to study the relation between task 1578 superposition and task composition. 1579

In Figure 8c (setting 4), we notice that the medians of the maximum non-task-answer probability are slightly higher than (or comparable to) those of the task answers for last\_letter and last\_letter\_cap. This could be due to the fact that these tasks are more challenging for the models. On the other hand, the medians of most probable non-task-answer probability are significantly lower than that of task answers for first\_letter and first\_letter\_cap, where we consider the model is performing task superposition on.

1586

1587 E.3 MEASURING ACCURACY IN TASK SUPERPOSITION

We further investigate how task superposition affect task performance. In particular, for a prompt consisting examples of K tasks (each task has an equal number of task examples), we define a task being correctly performed if its task answer lies in the top-K answers (that we use beam search to find). We compare the accuracy against individual task accuracy where we provide prompts consisting of task examples of only 1 task and define the task being correctly performed if the task answer is the top-1 answer.

We calculate accuracy in K = 1 case and K > 1 case using 100 prompts and show the result in Table 2. We find that as we increase the number of tasks, all models exhibit an accuracy decrease in correctly performing each individual task. Notably, however, the accuracy degradation of Llama3-70B is less than that of Llama2-70B across most tasks. We believe this is indicative that improved model training techniques lead to better preservation of task superposition,

Model		<i>K</i> =	= 1		K = 4					
	t1	t2	t3	t4	t1	t2	t3	t4		
GPT-3.5	100	99	85	91	95(-5)	90(-9)	84(-1)	84(-7)		
Llama-3 70B	100	99	97	99	100(0)	99(0)	96(-1)	92(-7)		
Llama-2 70B	100	96	69	77	88(-12)	96 (O)	63(-6)	46(-31)		
Qwen-1.5 72B	100	94	52	70	66(-34)	91(-3)	28(-24)	34(-36)		

(a) Setting 1: Addition in original numerical form and in different languages where t1 = add, t2 = add\_in\_en, t3 = add\_in\_fr, t4 = add\_in\_es.

Model		K = 1		K = 3				
	t1	t2	t3	t1	t2	t3		
GPT-3.5	100	100	100	93(-7)	62(-38)	56(-44)		
Llama-3 70B	100	100	100	82(-18)	90(-10)	83(-17)		
Llama-2 70B	97	94	90	75(-22)	75(-19)	50(-40)		
Qwen-1.5 72B	100	99	91	65(-35)	87(-12)	48(-43)		

(b) Setting 2: Naming the capital, continent and capitalize the country name where t1 = capital, t2 = continent, t3 = capitalization.

1646	Model			K = 1			 K = 3				
1647			_	1	-					-	
1648			t	51 1	t2	t3	tl	t2	t3		
1649	GPT-3	3.5	1	00 1	100	100	100(0)	97(-3)	97(-3)		
1650	Llama	-3 70B	1	00 1	100	100	99(-1)	100(0)	99(-1)		
1651	Llama	1570B	D 1	00 1	100	95 00	95(-5)	99(-1)	84(-11)		
1652	Qweii	-1.372	<b>D</b> 1	00 1	100	99	100 (0)	98 (-2)	98 (-1)	_	
1653	(c) Setti	ng 3: t	1 = c	opy (d	op1),	t2 =	copy (c	op2) and t3	= op1+op	2.	
1654											
1655	Model	K			$\zeta = 1$			K = 4			
1656		t1	t2	t3	t4		t1	t2	t3	t4	
1657	GPT-3.5	100	87	100	54	94	(-6)	56(-31)	97(-3)	12(-42)	
1658	Llama-3 70B	100	63	100	40	99	(-1)	29(-34)	99(-1)	13(-27)	
1659	Llama-2 70B	100	55	100	38	99	(-1)	35(-20)	97(-3)	7(-31)	
1660	Qwen-1.5 72B	100	62	100	34	89	(-11)	30(-32)	100(0)	15(-19)	

(d) Setting 4: First or last letter in upper or lower cases where t1 = first\_letter, t2 = last\_letter, t3 = first\_letter\_cap, t4 = last\_letter\_cap.

Table 2: Accuracy for each task in percentage, with the delta change given in parenthesis. For each setting we calculate the accuracy with prompts consisting of task examples of only one task (K = 1case) and with prompts consisting of examples from multiple tasks (K > 1 case).



Figure I: Accuracy for each choice of the intermediate layer  $\ell$  on task ret2 and plus2.



Figure II: Task vectors projected onto two axes chosen by LDA for two sets of tasks: (a) ret1,ret5 and ret7 and (b) plus1, plus5 and plus7.



Figure III: We vary the proportion,  $\lambda$ , between two tasks and observe how the output probabilities for the correct answers change. The proportion  $\lambda$  is varied in two ways: (1) in the top row, we plot the output from patching in a convex combination of task vectors for two tasks. (2) in the bottom row, we plot the output from a mixed proportion of in-context examples for the two tasks. Subplot (a) shows the output probabilities from mixing two retrieval tasks and (b) shows the probabilities from mixing two addition tasks.





Figure V: We vary the proportion,  $\lambda$ , between two tasks and observe how the output probabilities of a 1-head transformer for the correct answers change. The proportion  $\lambda$  is varied on two tasks: plus2 and plus5.