

# Test-Time Adaptation for LLMs via Hypernetwork-Generated LoRA Weights

Andrews George Varghese & Dhruv Kapur & Raj Maheshwari  
10-423/623 Generative AI Course Project

December 11, 2025

## 1 Abstract

We propose a novel approach for test-time adaptation of large language models (LLMs) using hypernetworks to predict Low-Rank Adaptation (LoRA) weights on-the-fly during inference. Unlike traditional fine-tuning methods that require separate adapters for each task, we propose a hypernetwork that generates specialized LoRA parameters for individual inputs in a single forward pass enabling dynamic, instance-level model specialization. Building on recent advances in text-to-LoRA generation which utilize a per-task label to generate the corresponding LoRA matrices, we aim to extend this to true per-input adaptation, potentially improving model performance across diverse benchmarks while maintaining computational efficiency. This would thus eliminate the need to create labels or summary prompts per task, signifying a shift from static model specialization per dataset to a more dynamic, prompt-conditional parameter generation.

## 20 1 Introduction

Large language models have achieved remarkable success across numerous natural language processing tasks, yet adapting these models to specialized domains or tasks typically requires expensive fine-tuning procedures. Recent work on Low-Rank Adaptation (LoRA) [Hu et al., 2022] has dramatically reduced the computational burden by constraining parameter updates to low-rank decompositions, reducing trainable parameters by factors as high as 10,000x while matching full fine-tuning performance.

However, even LoRA-based approaches face a fundamental limitation: each task requires a separate adapter that must be trained in advance. This creates significant engineering overhead when supporting multiple tasks and prevents dynamic adaptation to individual inputs. Recent advances in hypernetworks [Ha et al., 2017]—neural networks that generate weights for other networks—offer a compelling solution. Sakana AI’s Text-to-LoRA work [Charakorn

et al., 2025] demonstrates that hypernetworks can generate task-specific LoRA adapters from natural language descriptions, achieving 98% of specialized adapter performance without any task-specific training time.

We extend this paradigm by investigating **instance-level adaptation**: using hypernetworks to generate unique LoRA parameters for each input during inference, rather than uniform task-level adapters in order to maximize performance on every single input. This approach aims to combine the efficiency of hypernetwork-based weight generation with the fine-grained customization potential of per-input adaptation, potentially improving performance on complex tasks that exhibit high intra-task variance, while also studying the behavior of the generated LoRA weights on out of distribution tasks.

## 2 Dataset / Task

We will evaluate our approach on a diverse set of established benchmarks to assess both the quality and generalization capabilities of our hypernetwork-generated LoRA adapters:

**Reasoning Tasks:** ARC-Challenge Clark et al. [2018] (2,590 grade-school science questions requiring complex reasoning beyond fact retrieval), GSM8K Cobbe et al. [2021] (8,500 grade-school math word problems requiring 2–8 step solutions), and Winogrande Sakaguchi et al. [2020] (44,000 common-sense reasoning problems via pronoun resolution).

**Code Generation:** HumanEval Chen et al. [2021] (164 hand-written Python programming problems), and MBPP Austin et al. [2021] (974 basic Python tasks with assert-style test cases).

**Knowledge & Comprehension:** BoolQ Clark et al. [2019] (15,942 yes/no questions requiring inference from context), OpenBookQA Mihaylov et al. [2018] (5,957 science questions requiring multi-hop reasoning), and HellaSwag Zellers et al. [2019] (70,000 commonsense inference problems about physical scenar-

79 ios).

80 We use standard evaluation metrics for each bench-  
81 mark: accuracy for multiple-choice tasks (ARC,  
82 BoolQ, OpenBookQA, HellaSwag, WinoGrande),  
83 exact-match accuracy for GSM8K, and pass@1 (func-  
84 tional correctness) for code generation tasks.

85 As for model choice, we will begin experimenting  
86 with a pretrained model such as Pythia 70M/160M [Bi-  
87 derman et al., 2023] or SmoILM 135M [Allal et al.,  
88 2024], due their small size, and if time and compute  
89 budgets allow, we are also keen on experimenting with  
90 instruction-tuned models like SmoILM 135M Instruct  
91 or Gemma3 270M [Gemma, 2025].

### 92 3 Related Work

93 **Low-Rank Adaptation.** LoRA [Hu et al., 2022]  
94 established the currently most prevalent method of  
95 parameter-efficient fine-tuning by representing weight  
96 updates as  $\Delta W = BA$  where  $B \in \mathbb{R}^{d \times r}$  and  $A \in$   
97  $\mathbb{R}^{r \times k}$  with rank  $r \ll \min(d, k)$ . This reduces trainable  
98 parameters by 10,000x for GPT-3 175B while matching  
99 full fine-tuning performance.

100 **Hypernetworks.** The fundamental paradigm of us-  
101 ing one network to generate weights for another was  
102 established by Ha et al. [2017], demonstrating end-to-  
103 end training via backpropagation for generating non-  
104 shared LSTM weights. This concept was extended to  
105 transformers by HyperFormer Mahabadi et al. [2021],  
106 which generates adapter parameters for all layers by  
107 conditioning on task embeddings, layer position, and  
108 module type, achieving strong multi-task performance  
109 while adding only 0.29% parameters per task.

110 **Meta-Learning for Fast Adaptation.** Model-  
111 Agnostic Meta-Learning (MAML) [Finn et al., 2017]  
112 introduced gradient-based meta-learning through bi-  
113 level optimization, enabling rapid task adaptation with  
114 few gradient steps. While MAML requires gradient up-  
115 dates during adaptation, hypernetwork approaches can  
116 generate task-specific parameters in a single forward  
117 pass, representing complementary trade-offs between  
118 adaptation quality and computational efficiency.

119 **Text-to-LoRA and Test-Time Adaptation.** Most  
120 relevant to our work, Charakorn et al. Charakorn  
121 et al. [2025] showed that hypernetworks can gener-  
122 ate complete LoRA adapters directly from natural-  
123 language task descriptions, compressing hundreds of  
124 task-specific adapters into a single model. In parallel,  
125 HyperDecoders Ivison and Peters [2022] proposed a  
126 sequence-to-sequence architecture that decodes serial-  
127 ized parameters (e.g., adapters or classification heads)  
128 conditioned on natural-language instructions. Unlike  
129 fixed-format hypernetworks, HyperDecoders treat pa-  
130 rameter generation as an autoregressive text-generation

problem, which are then injected into a base model for  
131 task execution.  
132

While Text-to-LoRA operate at *task-level granularity*  
133 producing one set of parameters per task de-  
134 scription, HyperDecoder approaches only NLP sub-  
135 problems instead of language modeling tasks.  
136

Our work extends these ideas to *instance-level adapta-  
137 tion*, generating unique LoRA parameters for  
138 each inference input to maximize per-example perfor-  
139 mance on heterogeneous data in the language modeling  
140 paradigm.  
141

## 4 Approach

Our approach builds upon the Text-to-LoRA archi-  
143 tecture while introducing instance-level conditioning  
144 mechanisms:  
145

**Architecture.** We will implement a transformer-  
146 based hypernetwork that takes as input: (1) em-  
147 beddings of the input text using a frozen en-  
148 coder (e.g., gte-large-en-v1.5 [Li et al.,  
149 2023], ModernBERT-base [Warner et al., 2024]),  
150 NeoBERT [Breton et al., 2025], or alternatively a  
151 frozen copy of the target language model itself (2)  
152 learned module-specific embeddings identifying target  
153 layers (query and value projections), and (3) layer-  
154 specific positional information. (2) and (3) will build  
155 on work by Sakana [Charakorn et al. [2025]]. The hy-  
156 pernetwork processes these concatenated embeddings  
157 and generates LoRA weight matrices ( $B$  and  $A$ ) via  
158 linear projection heads.  
159

**Training Strategy.** We propose training the hy-  
160 pernetwork end-to-end on downstream tasks, back-  
161 propagating through the (frozen) language model and  
162 attached predicted adapter weights to the hypernet-  
163 work. This updates the hypernetwork’s parameters to  
164 learn meaningful mappings from input semantics to  
165 valid LoRA weight structures that maximize the per-  
166 formance of the language model on training examples.  
167 An alternate idea is to train the hypernetwork using a  
168 reconstruction loss, i.e. per-task LoRA adapters are  
169 learned using PEFT and saved to create a supervised  
170 dataset that maps inputs to their corresponding LoRA  
171 adapters. The hypernetwork is then trained to recon-  
172 struct LoRA weights from the input questions from  
173 the dataset. However, in additon to requiring hundreds  
174 LoRA adapters as training examples, [Charakorn et al.,  
175 2025] also shows that this method fails to generalize .  
176

**Instance-Level Conditioning.** Unlike Text-to-  
177 LoRA which conditions only on task descriptions, we  
178 will condition the hypernetwork on each individual in-  
179 put sequence. This enables the generation of input-  
180 specific LoRA weights that can adapt to intra-task var-  
181 iations, edge cases, and specific input characteristics,  
182

183 and crucially, doesn't require additional task labeling/description generation for the hypernetwork input.  
184

185 **Baseline Method.** Our primary baseline is standard  
186 LoRA fine-tuning on each benchmark dataset, which  
187 represents the current state-of-the-art for parameter-  
188 efficient adaptation, as well as testing for performance  
189 against the base model on out-of-distribution bench-  
190 marks.

191 **Key Contribution.** Our primary novel contribu-  
192 tion is extending hypernetwork-based LoRA genera-  
193 tion from task-level to *instance-level* adaptation, val-  
194 idating the feasibility of such test-time adaptation for  
195 LLMs, and identifying the feasibility of such a method  
196 for handling out-of-distribution inputs. This requires  
197 developing effective conditioning mechanisms that can  
198 extract input-specific features predictive of optimal  
199 LoRA adapter configurations, while tackling compu-  
200 tational efficiency problems such as handling different  
201 LoRA weights for each element in a batch during train-  
202 ing, as well as training issues such as avoiding overfit-  
203 ting to individual examples.

## 204 5 Expected Outcomes

205 We hypothesize that instance-level LoRA adaptation  
206 should have a reasonable performance in comparison  
207 to task-specific LoRA adapters, with the potential to  
208 outperform the baseline on benchmarks with high  
209 within-task diversity. Specifically:

210 **Performance Targets:** We expect to achieve 85-  
211 95% of standard LoRA performance (which itself  
212 matches or exceeds full fine-tuning) across bench-  
213 marks, with the potential to exceed task-level ap-  
214 proaches on complex reasoning tasks like Arc-  
215 Challenge and GSM8K where individual problems  
216 may benefit from specialized adaptations.

217 **Efficiency Gains:** Following Text-to-LoRA's re-  
218 sults, we anticipate computational efficiency compared  
219 to adapter training (just a single forward pass through  
220 the hypernetwork), while also potentially achieving  
221 better performance on benchmarks than in-context  
222 learning.

223 **Generalization:** We hope to demonstrate general-  
224 ization to held-out benchmarks, showing that the hy-  
225 pernetwork learns generalizable mappings from input  
226 characteristics to optimal LoRA configurations rather  
227 than memorizing task-specific solutions.

228 **Analysis:** We will analyze the learned latent repre-  
229 sentations to understand what characteristics of the in-  
230 puts drive LoRA parameter selection, visualize the di-  
231 versity of generated adapters across inputs, and identify  
232 which types of problems benefit most from instance-  
233 level adaptation through comparison with the baseline.

## 5.1 Potential Challenges

234 Key challenges include: (1) achieving stable training  
235 where the hypernetwork actually learns to predict valid  
236 and usable LoRA wieghts, (2) preventing overfitting  
237 to individual training examples, (3) balancing hyper-  
238 network capacity against memory constraints, (4) de-  
239 veloping robust regularization strategies and (5) max-  
240 imizing deliverables given our compute and time con-  
241 straints.  
242

## 243 6 Plan

244 There are 4 primary tasks that need to be completed for  
245 this project:

- 246 1. Building a flexible evaluation framework that we  
247 will use to evaluate the baseline language model,  
248 the task-specific fine-tuned models and our hyper-  
249 network augmented model on the selected bench-  
250 marks. Available options like Harness [Gao et al.,  
251 2024] and LightEval [Habib et al., 2023] provide  
252 reliable evaluation pipelines but are too opaque to  
253 be easily used in this project.
- 254 2. Creating the PEFT pipeline to train the per-task  
255 LoRA adapters.
- 256 3. Building and training the hypernetwork.
- 257 4. Analysing the generated LoRA adapters.

258 All team members will participate in hyperparameter  
259 tuning, debugging, and results interpretation. We will  
260 use pair programming for critical components like the  
261 hypernetwork forward pass and hypernetwork output  
262 analysis.

263 **Midway Executive Summary Deliverables:** By the  
264 midway point, we will have: (1) completed preliminary  
265 hypernetwork architecture implementation, (2) estab-  
266 lished evaluation pipeline and used it to evaluate the  
267 base model and (3) PEFT pipeline for per-benchmark  
268 fine-tuning.

## 269 References

270 Loubna Ben Allal, Anton Lozhkov, Elie Bakouch, Le-  
271 andro von Werra, and Thomas Wolf. Smollm - blaz-  
272 ingly fast and remarkably powerful, 2024.

273 Jacob Austin, Augustus Odena, Maxwell Nye, Maarten  
274 Bosma, Henryk Michalewski, Sharad Ramamurthy,  
275 David Bieber, Aleksandar Makelov, Chenjie Xiong,  
276 Michihiro Yasunaga, and Qian Lin. Program syn-  
277 thesis with large language models. *arXiv preprint arXiv:2108.07732*, 2021.

- 279 Stella Biderman, Hailey Schoelkopf, Quentin Gregory  
280 Anthony, Herbie Bradley, Kyle O'Brien, Eric Hal-  
281 lahan, Mohammad Aflah Khan, Shivanshu Purohit,  
282 USVSN Sai Prashanth, Edward Raff, et al. Pythia:  
283 A suite for analyzing large language models across  
284 training and scaling. In *International Conference on*  
285 *Machine Learning*, pages 2397–2430. PMLR, 2023.
- 286 Lola Le Breton, Quentin Fournier, Mariam El  
287 Mezouar, and Sarath Chandar. Neobert: A next-  
288 generation bert, 2025. URL <https://arxiv.org/abs/2502.19587>.
- 290 Rujikorn Charakorn, Edoardo Cetin, Yujin Tang,  
291 and Robert Tjarko Lange. Text-to-LoRA: Instant  
292 transformer adaption. In *Forty-second Interna-*  
293 *tional Conference on Machine Learning*, 2025.  
294 URL [https://openreview.net/](https://openreview.net/forum?id=zWskCdu3QA)  
295 [forum?id=zWskCdu3QA](https://openreview.net/forum?id=zWskCdu3QA).
- 296 Mark Chen, Jerry Tworek, Heewoo Jun, Qiming Yuan,  
297 Henrique de Oliveira Pinto, Jared Kaplan, Harri Ed-  
298 wards, Yuri Burda, Nicholas Joseph, Greg Brock-  
299 man, Kate Ray, Raahil Puri, Gretchen Krueger,  
300 Michael Petrov, Heidy Khlaaf, Girish Sastry, Pamela  
301 Mishkin, Brooke Chan, Scott Gray, Nick Ryder,  
302 Mikhail Pavlov, Ashish Power, Lukasz Kaiser, Mo-  
303 mhammad Bavarian, Robert Winter, Phil Tillet, Felipe  
304 Such, David Cummings, Matthias Plappert, Feyn-  
305 man Chantzis, Alex Barnes, Ariel Herbert-Voss,  
306 Will Guss, Kamal Dwivedi, Trent Millard, Shantanu  
307 Ogawa, Sam Thompson, Time Chen, Amos Azaria,  
308 Chris McLeavey, Edward Farhi, Josh Achiam, Sand-  
309 hini Agarwal, Andrew Lohn, Vicki Zhuang, Jack  
310 Clark, Dario Amodei, Ingrid Kaplan, Tom McCan-  
311 dlish, and Ilya Sutskever. Evaluating large lan-  
312 guage models trained on code. *arXiv preprint*  
313 *arXiv:2107.03374*, 2021.
- 314 Christopher Clark, Kenton Lee, Ming-Wei Chang, and  
315 Tom Kwiatkowski. Boolq: Exploring the surprising  
316 difficulty of natural yes/no questions. In *Proceed-  
317 ings of the 2019 Conference of the North American*  
318 *Chapter of the Association for Computational Lin-*  
319 *guistics*, 2019.
- 320 Peter Clark, Isaac Cowhey, Oren Etzioni, Tushar Khot,  
321 Ashish Sabharwal, Carissa Schoenick, and Oyvind  
322 Tafjord. Think you have solved question answering?  
try arc, the ai2 reasoning challenge. In *Proceed-  
323 ings of the 2018 Conference on Empirical Methods*  
324 *in Natural Language Processing*, 2018.
- 325 Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian,  
326 Mark Chen, Heewoo Jun, Lukasz Kaiser, Matthias  
327 Plappert, Jerry Tworek, Jacob Hilton, Reiichiro  
328 Nakano, Christopher Hesse, and John Schulman. Training verifiers to solve math word problems. *arXiv preprint arXiv:2110.14168*, 2021.
- 329 Chelsea Finn, Pieter Abbeel, and Sergey Levine. Model-agnostic meta-learning for fast adaptation of deep networks. In *Proceedings of the 34th International Conference on Machine Learning*, volume 70 of *Proceedings of Machine Learning Research*, pages 1126–1135. PMLR, 2017.
- 330 Leo Gao, Jonathan Tow, Baber Abbasi, Stella Biderman, Sid Black, Anthony DiPofi, Charles Foster, Laurence Golding, Jeffrey Hsu, Alain Le Noac'h, Haonan Li, Kyle McDonell, Niklas Muennighoff, Chris Ociepa, Jason Phang, Laria Reynolds, Hailey Schoelkopf, Aviya Skowron, Lintang Sutawika, Eric Tang, Anish Thite, Ben Wang, Kevin Wang, and Andy Zou. The language model evaluation harness, 07 2024. URL <https://zenodo.org/records/12608602>.
- 331 Team Gemma. Gemma 3. 2025. URL <https:// goo.gle/Gemma3Report>.
- 332 David Ha, Andrew M Dai, and Quoc V Le. Hyper-  
333 networks. In *International Conference on Learn-  
334 ing Representations*, 2017. URL <https:// openreview.net/ forum?id=rkpACe1lx>.
- 335 Nathan Habib, Clémentine Fourrier, Hynek Kydlíček,  
336 Thomas Wolf, and Lewis Tunstall. Lighteval: A  
337 lightweight framework for llm evaluation, 2023.  
URL <https://github.com/huggingface/lighteval>.
- 338 Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan  
339 Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang,  
340 and Weizhu Chen. LoRA: Low-rank adapta-  
341 tion of large language models. In *Internation-  
342 al Conference on Learning Representations*,  
343 2022. URL <https://openreview.net/ forum?id=nZeVKeefYf9>.
- 344 Hamish Ivison and Matthew E. Peters. Hyperde-  
345 coders: Instance-specific decoders for multi-task  
346 nlp, 2022. URL <https://arxiv.org/abs/2203.08304>.
- 347 Zehan Li, Xin Zhang, Yanzhao Zhang, Dingkun Long,  
348 Pengjun Xie, and Meishan Zhang. Towards general  
349 text embeddings with multi-stage contrastive learn-  
350 ing. *arXiv preprint arXiv:2308.03281*, 2023.
- 351 Rabeeh Mahabadi, Sebastian Ruder, Mostafa De-  
352 ghani, and James Henderson. Parameter-efficient  
353 multi-task fine-tuning for transformers via shared

377 hypernetworks. In *Proceedings of the 59th Annual*  
378 *Meeting of the Association for Computational Lin-*  
379 *guistics and the 11th International Joint Conference*  
380 *on Natural Language Processing (Volume 1: Long*  
381 *Papers)*, pages 565–576. Association for Compu-  
382 *tational Linguistics*, 2021.

383 Todor Mihaylov, Peter Clark, Tushar Khot, and Ashish  
384 Sabharwal. Can a suit of armor conduct electricity?  
385 a new dataset for open book question answering. In  
386 *Proceedings of the 2018 Conference on Empirical*  
387 *Methods in Natural Language Processing*, 2018.

388 Keisuke Sakaguchi, Ronan Le Bras, Chandra Bhaga-  
389 vatula, and Yejin Choi. Winogrande: An adversarial  
390 winograd schema challenge at scale. In *Proceedings*  
391 *of the AAAI Conference on Artificial Intelligence*,  
392 volume 34, pages 8732–8740, 2020.

393 Benjamin Warner, Antoine Chaffin, Benjamin Clavié,  
394 Orion Weller, Oskar Hallström, Said Taghadouini,  
395 Alexis Gallagher, Raja Biswas, Faisal Ladhak, Tom  
396 Aarsen, Nathan Cooper, Griffin Adams, Jeremy  
397 Howard, and Iacopo Poli. Smarter, better, faster,  
398 longer: A modern bidirectional encoder for fast,  
399 memory efficient, and long context finetuning and  
400 inference, 2024. URL <https://arxiv.org/abs/2412.13663>.

402 Rowan Zellers, Ari Holtzman, Yonatan Bisk, Ali  
403 Farhadi, and Yejin Choi. Hellaswag: Can a machine  
404 really finish your sentence? In *Proceedings of the*  
405 *57th Annual Meeting of the Association for Compu-*  
406 *tational Linguistics*, 2019.