NNSIGHT AND NDIF: DEMOCRATIZING ACCESS TO OPEN-WEIGHT FOUNDATION MODEL INTERNALS

Anonymous authors

Paper under double-blind review

ABSTRACT

We introduce NNsight and NDIF, technologies that work in tandem to enable scientific study of the representations and computations learned by very large neural networks. NNsight is an open-source system that extends PyTorch to introduce deferred remote execution. NDIF is a scalable inference service that executes NNsight requests, allowing users to share GPU resources and pretrained models. These technologies are enabled by the *intervention graph*, an architecture developed to decouple experiment design from model runtime. Together, this framework provides transparent and efficient access to the internals of deep neural networks such as very large language models (LLMs) without imposing the cost or complexity of hosting customized models individually. We conduct a quantitative survey of the machine learning literature that reveals a growing gap in the study of the internals of large-scale AI. We demonstrate the design and use of our framework to address this gap by enabling a range of research methods on huge models. Finally, we conduct benchmarks to compare performance with previous approaches. Code and documentation will be made available open-source.

004

010

011

012

013

014

015

016

017

018

019

021

1 INTRODUCTION

Research on large-scale AI currently faces two practical challenges: *lack of transparent model access*, and *lack of adequate computational resources*. Model access is limited by the secrecy of state-ofthe-art commercial model parameters (OpenAI et al., 2023; Anthropic, 2024; Gemini Team et al., 2024), and computational resources are limited by funding and engineering barriers. Commercial application programming interfaces (APIs) help amortize costs by offering frontier models "as a service". However, these APIs lack the transparency necessary to enable scientists to study model internals, e.g., by providing access to intermediate activations or gradients used during neural network inference and training.

Surveys of AI research needs have documented this situation (Shevlane, 2022; Casper et al., 2024).
 Bucknall & Trager (2023) highlight the need for *structured model access APIs* that offer greater
 transparency than existing commercial inference APIs. The current paper aims to meet this need by
 defining an expressive and efficient technical standard that addresses the key problems for providing
 a practical and useful structured model access API to support research on large models.

041 Our work makes three contributions:

The *intervention graph* architecture (Figure 1b), an approach for organizing experiments on very large models that reduces engineering burden, enhances reproducibility, and enables low-cost communication with remote models. We compare the intervention graph architecture to classical approaches for running similar experiments, and show how it achieves these goals.

NNsight (Figure 1a), an open-source implementation of the intervention graph architecture that extends PyTorch (Paszke et al., 2019) to create an expressive programming idiom that supports transparent model interventions on large-scale AI without requiring users to store or manage model parameters locally. NNsight builds intervention graphs using the same deferred computation graph idiom that deep learning frameworks adopt to enable automatic differentiation (Bottou & Le Cun, 1988; Bottou & Gallinari, 1990; Abadi et al., 2016; Al-Rfou et al., 2016; Bradbury et al., 2018).

NDIF (Figure 1c), an open-source cloud inference service that supports this idiom by providing behind-the-scenes user sharing of model instances to reduce the costs of large-scale AI research.



Figure 1: An example of the implementation of an NNsight intervention graph: (a) A user writes
 research code from which (b) an intervention graph is constructed. (c) The intervention operations are
 interleaved with the original model's computation and then executed. Values marked with .save()
 are made available to the user upon completion.

Crucially, by executing intervention graphs from multiple users with safe co-tenancy, NDIF supports the expressiveness of NNsight while avoiding costly startup times and communication costs incurred by traditional high-performance computing (HPC). We measure the performance of NDIF against the traditional approach, and we also compare it with Petals (Borzunov et al., 2023), a peer-to-peer approach to reducing barriers to distributed computing with large AI models.

2 SURVEYING MODEL AVAILABILITY AND RESEARCH USAGE

Compared to commercial APIs, models with openly downloadable parameters enable much more invasive experimentation and are a valuable resource for interpretability researchers (Scao et al., 2022; Biderman et al., 2023; Touvron et al., 2023; Jiang et al., 2023). To understand current research activity studying the internals of these open-weight models, we examine a set of 184 papers sourced from a recent survey of interpretability research (Ferrando et al., 2024) — see Appendix A for data curation details.

In Figure 2, we visualize the research usage of the largest open-weight models studied in these papers 088 over time. While there has been a rapid investment in *training* larger and more capable models, we 089 find that there is a disparity between the most capable systems¹ and those being *studied* in detail. 090 There is a small group of papers that study language models with $\geq 70\%$ MMLU performance (see (a) 091 in Figure 2), suggesting that this gap can partially be explained as a previous lack of sufficiently 092 capable models. However, even after the release of more-capable open-weight models, 60.6% of the surveyed research papers since February 2023 are still doing research on smaller, less performant 094 models (< 40% MMLU). This suggests that the gap is not just due to a lack of capable models, but 095 may also be due to engineering and infrastructural barriers.

096 Studying models with up to 70 billion parameters is possible for labs with access to high-end compute such as NVIDIA A100 or H100 80GB GPUs, but doing research beyond this scale represents a leap 098 in technical complexity and resource requirements. For example, inference on a 530 billion parameter 099 LLM uses 1TB of GPU memory, requiring many devices distributed across multiple nodes just to 100 load model parameters in 16-bit precision (Aminabadi et al., 2022; Dettmers et al., 2023). Even 101 if a scientist has access to such resources, they must integrate their experiment code with model 102 sharding, parallelization, quantization, and other deployment specifics. Thus, despite evidence that many fascinating capabilities appear only in the largest models (Brown et al., 2020; Rae et al., 2021; 103 Wei et al., 2022; Patel & Pavlick, 2022; Schaeffer et al., 2023), studying them can seem unattainable. 104

One solution is to study similarly capable, smaller models; this approach is used in two recent works that study Yi-34B (Young et al., 2024) and Qwen 2 72B (Yang et al., 2024) (see (b) in Figure 2) —

107

072 073

074

075

076

077

¹Defined in terms of MMLU performance (Hendrycks et al., 2021, see Appendix A for details.)



128 Figure 2: Most interpretability research is done on models that lag far behind the capabilities available in either closed- or open-models. Each blue point represents the MMLU performance of the largest open-weight 129 model studied by a surveyed paper, where the size of the point represents the model's parameter size. Models 130 without a recorded MMLU score were interpolated with nearest neighbors. There is a significant gap between 131 the performance of models studied (blue line) and the capabilities of leading open-weight models (shown in 132 orange). This gap is extended even further when considering the performance of leading closed-weight models 133 (black line). (a) A small group of papers study language models with ≥70% MMLU performance, which can account for at least some of the gap. However, many researchers are still studying smaller, less performant 134 models that hover around baseline performance (shown in gray). (b) Studying smaller, but still capable models 135 such as Qwen 72B or Yi-34B may be part of the solution to closing this research gap, but these models still 136 underperform the leading open-weight model, Llama 3.1 405B. 137

these are models that use distillation techniques to compress larger model capabilities into smaller
model sizes. However, these models still underperform the most capable open-weight models, such
as Llama 3.1 405B (Dubey et al., 2024). In the following sections we describe how NNsight and
NDIF are designed to address this research gap by providing a flexible system for studying very large
open models, presenting a reusable and modular service suitable for fully-transparent large-scale
deployment distributed across many devices and nodes.

- 145
- 146 147

3 A FRAMEWORK FOR EXPERIMENTS ON LARGE-SCALE AI

Research on the internals of pretrained neural networks follows one common pattern: Experimental code is inserted between computational steps of the network, interleaved with code previously written to run the network. Depending on the experiment, researchers' code may inspect internal representations (Bau et al., 2017), perform interventions (Vig et al., 2020), collect gradients (Sundararajan et al., 2017; Selvaraju et al., 2017; Ancona et al., 2019; Bastings et al., 2022; Kramár et al., 2024), modify parameters (Meng et al., 2022), or use and train the parameters of supplementary models (Lester et al., 2021; Wu et al., 2023; Huben et al., 2024; Templeton et al., 2024; Marks et al., 2024).

Traditionally, researchers have organized their code by creating bespoke versions of neural network modifications and hooks for each experiment using callback systems (Li et al., 2024; Wang et al., 2022; Geva et al., 2023, and see Figure 4a), or by designing custom neural network systems with built-in access points (Nanda & Bloom, 2022; Wu et al., 2024b). While tying experiments with network implementations is primarily a matter of code organization for small models, it burdens researchers with the responsibility and cost of model deployment and storage, increasing engineering challenges as parameter sizes increase. This pattern also makes it difficult for external research groups to reproduce large distributed neural network experiments without duplicating the environment.



Figure 3: Augmenting the computation graph C (pink graph on the left) with an intervention graph (green, right). (a) A getter edge brings the output of a variable node v_1 in the computation graph to an apply node a'_1 in the intervention graph. (b) A setter edge sends the output of v'_2 in the intervention graph back to a_3 in the computation graph.

188

189

190

191

192

193 194

196

197

199

200

201

We introduce a new framework that implements the following pattern: Separate experimental and engineering code, store the experiment, and then execute it on shared computational infrastructure. We achieve this separation by expressing an experiment as an *intervention graph* that can be saved, transmitted, optimized, and interleaved with model execution.

This framework implements three core requirements necessary for scalable, reproducible researchthat are not currently met by classical approaches:

- **Separation**: The framework must decouple experimental and engineering code, so that the increasingly complex challenge of running the underlying model can be tackled separately from experimental design.
- **Expressiveness**: The framework must allow the user to express all the experiments that they would have been able to express if they were using PyTorch on local hardware.
 - **Co-Tenancy**: Experiments must be able to be safely, efficiently, and concurrently combined with model execution at runtime, to allow the service to amortize costs across many users.

Each of these requirements is addressed by one of the aspects of our framework: First, the *intervention graph* allows us to decouple research and engineering code. Next, the *NNsight API* is designed to be fully expressive, because users define interventions using Python and PyTorch code. Finally, the *NDIF inference server* enables co-tenancy, by setting up a central system to which experiments can be safely transmitted, optimized, and run concurrently.

²⁰² 203 3.

3.1 Representing experiments as graphs

Our core innovation is a portable, serializable representation of an experiment, called the *intervention graph*, which dynamically captures modifications to the model's computation graph. This graph can be stored in JSON format, version-controlled, optimized, and sent to or retrieved from remote systems for execution. In this section, we formalize the notion of an intervention graph.

Computation Graph. Our formalism adopts Al-Rfou et al. (2016, Theano)'s definition of a computation graph². The computation graph is a weakly connected, bipartite, directed, acyclic graph defined as C = (V, A, E), where V and A are variable nodes and apply nodes, respectively. E is a set of directed edges connecting the variable nodes and the apply nodes. The definition and properties of V, A, and E are as follows:

- 214 215
- $V = \{v_1, \dots, v_m\}$ is the set of variable nodes representing objects. They are one-to-many, and are connected to apply nodes.

• $A = \{a_1, \dots, a_n\}$ is the set of apply nodes representing operations on variable nodes. They are many-to-one, and are connected to variable nodes.

218 219

216

217

220

221

222

223

224 225

226

227 228

229

230

231 232

233

234 235

236

245

• $E = \{e_1, \dots, e_p\}$ is the set of directed edges. An edge e in E either connects a variable node v to an apply node a, or vice versa. Thus, $E \subseteq (V \times A) \cup (A \times V)$ — making C a bipartite graph as there are no connections between nodes of the same type.

In our context, a node that is "one-to-many" has a single input edge, but multiple output edges, as seen with variable nodes. Likewise, a node that is "many-to-one" has multiple input edges but only a single output edge, as with apply nodes. In NNsight, apply nodes correspond to PyTorch Modules.

The graph C is *weakly connected* since converting all directed edges to undirected edges results in a connected graph. It is *bipartite* because edges only exist between variable and apply nodes.

Intervention Graph. Interventions are modifications to the computation graph C that introduce additional nodes and edges. We define an intervention as $I = (V', A', E', \mathcal{G}, \mathcal{S})$, where (V', A', E') forms an *intervention component* G. Each G inherits the properties of a computation graph. The sets \mathcal{G} and \mathcal{S} , referred to as getters and setters, define how G connects to C.

Integrating an intervention component into C is called *interleaving*, and consists of two parts:

- The getters (\mathcal{G}) connect variable nodes in C to apply nodes in G, i.e., $\mathcal{G} \subseteq V \times A'$.
- The setters (S) connect variable nodes in G back to apply nodes in C, i.e., $S \subseteq V' \times A$.

The *intervention graph* of NNsight is the union of all individual intervention components, $\mathcal{I} = \bigcup_{1 \le i \le k} G_i$. The final augmented computation graph C' is constructed by interleaving each intervention component with the original graph. C' is a new computation graph, inheriting the properties of C. When an intervention graph is sent to NDIF, the getters and setters are also included. Appendix B discusses implementation details.

For the intervention graph to be valid, for every *getter* edge $(v_i, a'_j) \in \mathcal{G}$ and *setter* edge $(v'_k, a_l) \in \mathcal{S}$, there must be no directed path from a_l to v_i or from v'_k to a'_j . This rule ensures that the augmented computation graph is acyclic.

246 3.2 FAMILIAR AND EXPRESSIVE INTERVENTIONS

To define and execute intervention graphs in code, the NNsight API uses PyTorch notation inside of a
Python context window. We enable the tracing of user-defined experiments inside the context window
by overriding basic Python operations as well as wrapping every PyTorch module and submodule to
expose their inputs, outputs, and gradients. We call this the *tracing context*.

In Figure 4, we compare two code snippets that define the same intervention — one using standard PyTorch hooks (Figure 4a) and the other using NNsight's API (Figure 4b). In both snippets, a language model (Llama-3.1-8B; Dubey et al., 2024) is loaded from HuggingFace (Wolf et al., 2020) and an intervention is specified – activating three neurons which cause the model to invert its output (e.g., producing "lie" instead of "truth"). While using standard PyTorch hooks requires creating custom hooks for each access point (Figure 4a, lines 8-13), NNsight's trace context window provides access to the intermediate inputs and outputs of all PyTorch modules (Figure 4b, lines 7-9). See Appendix C for additional NNsight code examples.

The experiment's corresponding intervention graph is defined upon exiting the context window, and execution does not occur until the context is complete. The graph can then either be executed locally, or it can be serialized and sent to a remote system for execution. Tensors within the deferred execution block are not directly available to the main code, but every deferred tensor provides a . save() method to make its data available to the main program. Although the tracing context defers execution, it is still possible to debug many issues locally by using the PyTorch FakeTensor system, which precomputes and checks tensor shapes and datatypes while building the computation graph.

Because PyTorch code is used inside the tracing context, any experiment that can be written using standard PyTorch hooks can also be expressed using NNsight and will feel familiar to PyTorch users.

 $^{^{2}}$ Al-Rfou et al. (2016) allow apply nodes to have many outgoing edges, whereas in ours, apply nodes can only have one outgoing edge. Appendix E argues that the two definitions are equivalent for our use case.

```
270
            1 from transformers import AutoTokenizer, AutoModelForCausalLM
                                                                                1 from nnsight import LanguageModel
              model_id = "meta-llama/Meta-Llama-3.1-8B"
tokenizer = AutoTokenizer.from_pretrained(model_id)
                                                                                    model_id = "meta-llama/Meta-Llama-3.1-8B"
271
                                                                                 3 lm = LanguageModel(model_id)
272
           4 lm = AutoModelForCausalLM.from_pretrained(model_id)
273
            5 mlp = lm.model.layers[16].mlp.down_proj
                                                                                 4 mlp = lm.model.layers[16].mlp.down_proj
                                                                                 5 neurons = [394, 5490, 8929]
6 prompt = "The truth is the"
              neurons = [394, 5490, 8929]
prompt = "The truth is the"
274
275
276
            8 def pre_hook_fn(module, input):
277
           9
                   input[0][:,-1,neurons] = 10
278
           10 hook = mlp.register_forward_pre_hook(pre_hook_fn)
                                                                                 7
                                                                                    with lm.trace(prompt, remote=True):
           11
             inputs = tokenizer(prompt, return_tensors="pt")
                                                                                        mlp.input[:, -1, neurons] = 10
279
           12 out = lm(**inputs)
                                                                                 9
                                                                                        out = lm.output.save()
280
           13 hook.remove()
281
           14 last = out["logits"][:, -1].argmax()
                                                                                10 last = out["logits"][:, -1].argmax()
              prediction = tokenizer.decode(last)
                                                                                11 prediction = lm.tokenizer.decode(last)
282
           15
                                                                                12 print(prediction)
           16 print(prediction)
283
                                             (a)
                                                                                                            (b)
284
           Figure 4: Experiment code expressed using (a) standard PyTorch hooks and (b) the NNsight API.
285
```

Figure 4: Experiment code expressed using (a) standard PyTorch hooks and (b) the NNsight API.
Both code snippets define the same intervention – activating three neurons which cause the model to invert the meaning of its output (e.g., producing "lie" rather than "truth"). The PyTorch intervention code captured in six lines of code (a, lines 8-13), can be easily expressed using NNsight with three lines of code (b, lines 7-9). Standard PyTorch requires creating custom hooks for each access point, whereas with NNsight, all module inputs and outputs can be accessed within a single trace context.

291

298

299

300

NNsight wraps all 217 fundamental PyTorch tensor operations, supports built-in modules and custom
 neural network architectures, and includes infrastructure for loading models distributed by popular
 platforms such as HuggingFace. Users can express a wide range of interventions beyond simple
 probing and manipulation, including custom attention mechanisms, dynamic computation graph
 alterations, weight updates, and training loops, as well as computing gradients on backwards passes
 and gradient flow modifications.

Code Example 1 shows an NNsight implementation for attribution patching (Kramár et al., 2024), a method that involves both causal interventions and computation of gradients within the model.

```
get_batch_states_and_grads(model, submodules, submod_names, batch, batch_target_ids):
         1
           def
301
                hidden_states_cache = {}
         3
                grads_cache = {}
302
303
         5
                with model.trace(batch, remote=True) as tracer:
                   for submodule in submodules:
    if 'neurons' in submod_names[submodule]:
304
         6
305
                            x = submodule.input
         8
                        else:
         9
306
                            x = submodule.output[0]
         10
307
                        hidden_states_cache[submodule] = x.save()
                        grads_cache[submodule] = x.grad.save()
308
                   neg_log_likelihood = metric_nll(model, batch_target_ids).save()
309
         14
                    neg_log_likelihood.sum().backward()
310
         16
                hidden_states = {k : v.value[:,token_idx] for k, v in hidden_states_cache.items()}
311
                grads = {k : v.value[:,token_idx] for k, v in grads_cache.items()}
         17
                return hidden_states, grads, neg_log_likelihood
312
         18
```

Code Example 1: An NNsight implementation for attribution patching (Kramár et al., 2024)

313 314 315

316

317

318

Crucially, by separating experiment code from the underlying neural network implementation, NNsight allows remote execution of experiments by simply adding a remote=True flag, as in Code Example 2. This sends the experiment to NDIF for execution.

```
319 1 with lm.trace(prompt, remote=True):
320 3
321 3 out = lm.output.save()
```

Code Example 2: Setting remote=True in a tracing context passes the intervention graph to NDIF infrastructure to be interleaved remotely. This example assumes all objects are the same as in Figure 4.

327

331 332

333

334

335

336

337



Figure 5: Schematic of research community use of NDIF vs. HPC and Petals. Green nodes show custom experiments. NDIF (left) allows many researchers to share a common inference service that runs customized experiments with shared memory. In HPC (center), researchers are responsible for weight-loading and handling model memory overhead on their own separate instances. In peerto-peer swarm approaches like Petals (right), while GPU resources are shared, hidden states must be transferred between nodes during inference and returned to the user for custom interventions, resulting in costly data transfers.

338 339

340 3.3 **CO-TENANCY AND REMOTE INFRASTRUCTURE** 341

342 The NDIF service, in contrast to the NNsight system, is a multi-user runtime infrastructure designed 343 to amortize costs and remove engineering burden from scientists by concurrently sharing inference. 344 NDIF accepts intervention graphs created by NNsight, and executes them on large, preloaded models. 345

Unlike traditional distributed computing approaches for supporting research experiments in AI, the 346 NDIF infrastructure is designed to minimize the communication overhead involved in remote model 347 execution. Figure 5 compares the infrastructure to traditional approaches. 348

349 In high-performance computing services (HPC), shared computational resources are allocated to 350 researchers at the level of machines or virtual machines. That means that when researchers conduct 351 experiments on customized AI models, all of the engineering code, weight-loading, and model storage must be handled by the researcher on instances that are exclusively allocated to them. 352

353 One approach that has been proposed for sharing costly computational resources while retaining 354 researcher control is to distribute work using a peer-to-peer distributed service. This is the approach 355 supported by Petals (Borzunov et al., 2023), which provides a swarm of inference servers that 356 preload and provide inference computation services for layers of large neural networks. The Petals architecture can support custom experiment work by enabling researchers to host their own customized 357 intervention code for particular layers on their own nodes (e.g., on their own workstation) while 358 relying on the swarm for other inference steps. However, this architecture incurs costly data transfers 359 multiple times during a forward pass for experiments that intervene on the inference process. 360

361 NDIF is designed to execute intervention graphs within the high-performance cluster itself. Models 362 and engineering code are preloaded onto the cluster, and large model-internal data transfers are 363 avoided since custom experiment code does not need to be run on the researchers' own hardware. Critically, unlike traditional HPC and peer-to-peer approaches, NDIF is designed to allow multiple re-364 searchers to conduct customized experiments while sharing the underlying preloaded model instances on shared computing hardware. 366

367 Safe co-tenancy. NDIF assumes responsibility for ensuring that user experiments conducted on 368 the same model and hardware do not interfere with each other by preventing direct access to model 369 parameters and by enforcing a server-side whitelist of permitted operations that ensure complete separation of user experiments. The computation graph's interleaving mechanism allows users to 370 simulate parameter interventions without accessing or modifying the parameters themselves. NDIF 371 also tracks and authorizes individual usage to maintain fair resource allocation. 372

373 Compute efficiency. NDIF hosts a selection of preloaded models. For each model, a single shared 374 instance is open to minimize performance difficulties associated with weight-loading and model 375 startup. Researchers can access these models and conduct experiments on them through intervention requests conducted with the NNsight API. The infrastructure implements horizontal scaling and 376 dynamic resource allocation to support multiple user requests and ensure the model deployments 377 are operational and healthy at any given time. The NDIF server facilitates dynamic, bi-directional

communication with the client side, enabling users to retrieve and save requested results, and allowing
 users to minimize additional overhead when running experiments remotely compared to locally. This
 approach reduces costs for individual researchers, as they no longer need to set up and maintain their
 own separate high-performance computing environments. Figure 6 compares NDIF performance to
 using HPC and Petals.

The code for creating NDIF infrastructure will be made publicly available on Github, allowing users to create their own NDIF infrastructure for specialized use cases.

4 PERFORMANCE AND EVALUATION

In this section, we evaluate the performance of NNsight when used remotely, both using exclusive allocations in the traditional high-performance computing setting, and using the shared NDIF infer-ence service. We also compare to the Petals (Borzunov et al., 2023) distributed inference system, which can also reduce the costs of customized inference on huge models. Our evaluation focuses on the time required to load the model weights into memory, as well as the runtime of activation patching, a standard model intervention technique (Vig et al., 2020). For this purpose, we use a single batch of 32 examples from the Indirect Object Identification (IOI) dataset (Wang et al., 2022). We evaluate performance in three different settings: First, we measure the runtime when the models and interventions are executed entirely on a high-performance computing node (HPC). Then, we execute the models and interventions on the NDIF infrastructure, and contrast it against HPC execution times. Finally, we compare NDIF with Petals (Borzunov et al., 2023), and evaluate their relative performance on standard inference tasks as well as on intervention tasks.

High-Performance Computing. We first evaluate NNsight in HPC scenarios by comparing it against other popular libraries designed for intervening on the internal states of PyTorch models, namely TransformerLens (Nanda & Bloom, 2022), pyvene (Wu et al., 2024b), and baukit (Bau, 2022). Our comparison spans models ranging from 1.5 to 8.5 billion parameters. We consider GPT2-XL (Radford et al., 2019), Gemma-7B (Team et al., 2024), and Llama-3.1 8B (AI@Meta, 2024). The experiments were conducted on a single HPC node with four NVIDIA H100 PCIe 82GB GPUs with CUDA version 12.3 and an Intel(R) Xeon(R) Gold 6342 CPU with 24 cores.

The results are presented in Table 1. We find that the other libraries achieve comparable weightloading and activation patching times³ to NNsight run without NDIF, for all three models. Therefore
we argue that when using exclusive HPC allocations, the choice of library should primarily be guided
by factors such as usability, feature set, and compatibility with existing workflows.

Table 1: Runtime comparison of baukit, pyvene, and TransformerLens with NNsight for loading
PyTorch models into memory and intervening on their internal states. NNsight achieves comparable
performance with other frameworks in both evaluation settings across all three models.

Framework	Setup Time			Activation Patching		
	GPT2-XL	Gemma-7B	Llama-3.1-8B	GPT2-XL	Gemma-7B	Llama-3.1-8B
baukit	3.146 ± 0.006	6.112 ± 0.074	4.529 ± 0.049	0.073 ± 0.001	0.220 ± 0.006	0.352 ± 0.006
oyvene	3.143 ± 0.011	6.263 ± 0.487	5.736 ± 0.717	0.074 ± 0.001	0.222 ± 0.001	0.348 ± 0.006
TransformerLens	8.991 ± 0.024	21.560 ± 0.385	20.625 ± 0.365	0.205 ± 0.027	0.208 ± 0.002	0.332 ± 0.002
NNsight	3.532 ± 0.172	6.351 ± 0.062	5.991 ± 0.192	0.073 ± 0.002	0.227 ± 0.009	0.335 ± 0.010

Remote Execution. We measure the same evaluation settings for HPC versus remotely on an NDIF server. We compare the two setups using the OPT suite of language models (Zhang et al., 2022), ranging from 125 million to 66 billion parameters. For the experiments, we use the same infrastructure as in the comparison against other libraries. The results of this analysis are presented in Figure 6a & 6b. As model size increases, the time required to load the model into memory grows nearly linearly on HPC. In contrast, the time for loading is negligible when using NDIF, since models are preloaded by the service. We also observe that the remote execution on NDIF introduces a roughly constant overhead from communication between the local device and NDIF during activation

 ³There is one exception, as loading the model weights into memory takes about three times as long with
 TransformerLens as with other libraries. This is likely a result of the preprocessing steps to convert weights into a standardized format across different models.



(a) **HPC vs. NDIF.** We compare the 443 time required to set up the model 444 when NNsight is executed on an HPC node and an NDIF server. For 445 HPC execution, the model must be 446 loaded, so the setup time scales 447 nearly linearly with model size. 448 In contrast, remote execution with 449 NDIF bypasses loading weights.

(b) HPC vs. NDIF. We also com- (c) Petals vs. NDIF. We compare the runtime of activation patching when NNsight is executed on an HPC node, versus using a NNsight locally to remotely access an NDIF server. We observe that the remote execution introduces a roughly constant overhead from communication Petals for remote interventions. between the local device and NDIF.



Figure 6: Performance evaluation of NNsight and NDIF. In all experiments, the sample size is n = 128. Plot marker widths are comparable to the 95% confidence intervals.

452 453

450

451

454 patching, regardless of model size (see Table 2 in Appendix D). In our setting, remote execution 455 provides runtime improvements for models with 3 billion parameters and more, indicating that remote 456 execution via NDIF becomes increasingly beneficial as parameter size grows. 457

Distributed inference. We also compare NDIF against Petals, an open-source service for distributed 458 model inference (Borzunov et al., 2023). Petals enables inference on large language models across 459 shared resources by transmitting intermediate hidden states. To ensure a fair comparison, we deployed 460 private instances of both Petals and NDIF on a server with a single NVIDIA RTX A6000 49 GB 461 GPUs with CUDA version 12.5 and an AMD Ryzen 9 5900X CPU with 12 cores. The communication 462 between these instances took place via a network with a bandwidth of about 60 MB/s. As Petals was 463 primarily designed for standard remote inference, we first compare the two services in this context. 464 With Petals, the client transmits token embeddings to the server and receives final hidden states, 465 which can then be mapped to token probabilities. In contrast, NDIF sends the input data and the 466 intervention graph to the server. While NDIF could simply return the next token prediction, we opted to return the final hidden states to the client for a fair comparison. The results are presented in 467 Figure 6c. In this scenario, both frameworks demonstrate comparable performance. 468

469 Petals also support certain types of inference-time interventions on model internals. Therefore, we 470 also compare NDIF and Petals in executing activation patching. Petals allows returning intermediate 471 hidden states, which can be used for interventions on the model's residual stream. However, as Petals does not support server-side interventions, modifications to the internal state must be performed locally 472 on the client device. This process involves receiving hidden states at a specific layer, performing 473 local modifications, and then sending the modified hidden states back to the server. In contrast, NDIF 474 supports server-side interventions and metric computations. This allows us to avoid transmitting 475 hidden states between client and server by computing patching metrics (e.g., logit difference) on the 476 server and returning only those. The results are presented in Figure 6c. In this scenario, we find that 477 NDIF significantly outperforms Petals. 478

479 480

481

5 **RELATED WORK**

482 Research model hosting frameworks. The most similar previous works are Garçon (Elhage et al., 483 2021) and Petals (Borzunov et al., 2023). Garçon is a proprietary internal research system at Anthropic that supports remote inspection and customization of large models hosted on a remote cluster. Unlike 484 the computation-graph approach of NNsight, Garçon operates by allowing researchers to hook models 485 via arbitrary-code execution. As a result, co-tenancy is unsafe, so unlike the NDIF server that shares

model instances between many users, Garçon requires each user to allocate their own model instance,
 which is computationally expensive. The Petals system has similar goals, but instead of distributing
 experiment code, it distributes foundation-model computation across user-contributed nodes in a
 BitTorrent-style swarm. Petals achieves co-tenancy by leaving researcher-defined computations on
 the user's own system. Unlike Petals, NNsight can avoid costly model-internal data transfers by
 executing computation graphs on NDIF servers.

492 Several other solutions for sharing hosted model resources have been developed, including VLLM 493 (Kwon et al., 2023), Huggingface TGI (Dehaene et al., 2024) and Nvidia TensorRT-LLM (Nvidia, 494 2024). These are systems for LLM inference acceleration that support large-scale multiuser model 495 sharing. However, unlike NDIF, these systems only provide black-box API access to models, and do 496 not permit model inspection or modification, which limits their utility for interpretability research. Also related to our work are systems such as S-LoRA (Sheng et al., 2024), dLoRA (Wu et al., 2024a), 497 and Punica (Chen et al., 2024), which excel at efficient serving of many pretrained LoRA adapters in 498 parallel. While these LoRA-focused systems enable one form of model modification, NNsight and 499 NDIF provide finer-grained control that allows for a wider range of model interventions. 500

501 Frameworks for model internals. Numerous efforts have been made to create robust tools for 502 exploring and manipulating model internals. These tools are all designed to support *local* experiments, 503 executed on compute resources controlled by the researcher conducting them. TransformerLens (Nanda & Bloom, 2022) is designed to facilitate the exploration of GPT-style decoder-only language 504 models. It allows users to apply custom functions to outputs of specific model modules. Unlike 505 NNsight, which operates on arbitrary PyTorch networks, TransformerLens is limited to transformer-506 based language models. Pyvene (Wu et al., 2024b) supports configurable intervention schemes on 507 PyTorch modules, decorating a PyTorch module with hooks that allow activations to be collected 508 and overwritten. Pyvene operates at a higher level of abstraction, and can be configured as a 509 layer over NNsight. Baukit (Bau, 2022) provides a range of utilities that simplify tracing and 510 editing activations for local models. However, unlike NNsight's intervention graph, it lacks an 511 intermediate representation that would allow for interaction with larger, remotely-hosted models. 512 Penzai (Johnson, 2024) is an open-source functional library which provides a modular system that 513 allows for introspection, visualization, and manipulation of neural network computation graphs; Penzai works within the JAX framework, in contrast to NNsight which works with PyTorch models. 514

515 516

517

6 DISCUSSION

518 The growing scale of highly-capable AI systems has presented our field with an access challenge: 519 The most capable systems are increasingly costly to deploy, and that has meant that many of the most 520 interesting and important research targets have become difficult for scientists to study in depth. In 521 this paper we have proposed an architecture for reducing barriers to research by defining a technical 522 approach for separating experiment definition from network deployment. Together, the intervention 523 graph architecture, the NNsight API, and the NDIF service enable researchers to perform complex experiments and explore neural network internals at a scale previously limited by computational, 524 engineering, and financial constraints. 525

526 While our approach addresses technical issues for conducting research on very large open-weight 527 models, a key limitation is that it does not provide access to closed-parameter, proprietary models 528 like GPT-4 or Claude. However, proprietary vendors often have special requirements. For instance, they may wish to maintain secrecy and security of model parameters for business or safety reasons. 529 Our system is designed with these concerns in mind, and it is technically feasible for companies to 530 support under these requirements. We encourage commercial providers to consider adopting our 531 approach as part of their products to enable transparency and facilitate future scientific progress at 532 the frontiers of AI. 533

NDIF and NNsight define a way of doing research that acknowledges the important role of community
resources. Scientists in other fields have opened up new research vistas by pooling investments in
shared research instruments, databases, and facilities. By enabling the AI research community to
similarly pool its resources for large-scale inference, we can open a pathway for conducting research
on AI at scales that would be unattainable for individual research groups.

540 **ETHICS STATEMENT** 7 541

Our work is intended to democratize access to large neural networks by providing tools that make it easier for researchers to audit, understand, and interpret large AI models. We caution that such transparency may also enable abuse of large models, for example exploiting model vulnerabilities. Therefore, we strongly urge the responsible and ethical use, deployment, and monitoring of our tools, while emphasizing the importance of continued research to improve model interpretation.

546 547 548

549

551

554 555

556

542

543

544

545

Reproducibility Statement 8

550 To promote transparency and ensure the reproducibility of our results, the complete NNsight and NDIF frameworks, along with comprehensive documentation and all relevant experiment code, will 552 be made publicly available as open-source resources. This will allow researchers and practitioners to 553 replicate our framework, related experiments and build upon our work.

References

- Martín Abadi, Paul Barham, Jianmin Chen, Zhifeng Chen, Andy Davis, Jeffrey Dean, Matthieu Devin, 558 Sanjay Ghemawat, Geoffrey Irving, Michael Isard, et al. TensorFlow: a system for Large-Scale 559 machine learning. In 12th USENIX symposium on operating systems design and implementation (OSDI 16), pp. 265–283, 2016. 560
- 561 AI@Meta. Llama 3 model card, 2024. URL https://github.com/meta-llama/llama3/blob/ 562 main/MODEL_CARD.md. 563

Rami Al-Rfou, Guillaume Alain, Amjad Almahairi, Christof Angermüller, Dzmitry Bahdanau, Nicolas Ballas, Frédéric Bastien, Justin Bayer, Anatoly Belikov, Alexander Belopolsky, Yoshua Bengio, 565 Arnaud Bergeron, James Bergstra, Valentin Bisson, Josh Bleecher Snyder, Nicolas Bouchard, 566 Nicolas Boulanger-Lewandowski, Xavier Bouthillier, Alexandre de Brébisson, Olivier Breuleux, 567 Pierre Luc Carrier, Kyunghyun Cho, Jan Chorowski, Paul Francis Christiano, Tim Cooijmans, 568 Marc-Alexandre Côté, Myriam Côté, Aaron C. Courville, Yann Dauphin, Olivier Delalleau, Julien 569 Demouth, Guillaume Desjardins, Sander Dieleman, Laurent Dinh, Mélanie Ducoffe, Vincent Du-570 moulin, Samira Ebrahimi Kahou, D. Erhan, Ziye Fan, Orhan Firat, Mathieu Germain, Xavier Glorot, 571 Ian J. Goodfellow, M. Graham, Çaglar Gülçehre, Philippe Hamel, Iban Harlouchet, Jean-Philippe 572 Heng, Balázs Hidasi, Sina Honari, Arjun Jain, Sébastien Jean, Kai Jia, Mikhail Korobov, Vivek 573 Kulkarni, Alex Lamb, Pascal Lamblin, Eric Larsen, César Laurent, Sea Sun Lee, Simon Lefrançois, 574 Simon Lemieux, Nicholas Léonard, Zhouhan Lin, Jesse A. Livezey, Cory Lorenz, Jeremiah Lowin, 575 Qianli Ma, Pierre-Antoine Manzagol, Olivier Mastropietro, Robert T. McGibbon, Roland Memisevic, Bart van Merrienboer, Vincent Michalski, Mehdi Mirza, Alberto Orlandi, Christopher Joseph 576 Pal, Razvan Pascanu, Mohammad Pezeshki, Colin Raffel, Daniel Renshaw, Matthew Rocklin, 577 Adriana Romero, Markus Roth, Peter Sadowski, John Salvatier, François Savard, Jan Schlüter, John 578 Schulman, Gabriel Schwartz, Iulian Serban, Dmitriy Serdyuk, Samira Shabanian, Etienne Simon, 579 Sigurd Spieckermann, S. Ramana Subramanyam, Jakub Sygnowski, Jérémie Tanguay, Gijs van 580 Tulder, Joseph P. Turian, Sebastian Urban, Pascal Vincent, Francesco Visin, Harm de Vries, David 581 Warde-Farley, Dustin J. Webb, Matthew Willson, Kelvin Xu, Lijun Xue, Li Yao, Saizheng Zhang, 582 and Ying Zhang. Theano: A python framework for fast computation of mathematical expressions. 583 ArXiv, abs/1605.02688, 2016. URL https://api.semanticscholar.org/CorpusID:8993325. 584

- Reza Yazdani Aminabadi, Samyam Rajbhandari, Minjia Zhang, Ammar Ahmad Awan, Cheng Li, 585 Du Li, Elton Zheng, Jeff Rasley, Shaden Smith, Olatunji Ruwase, and Yuxiong He. Deepspeed-586 inference: Enabling efficient inference of transformer models at unprecedented scale. SC22: International Conference for High Performance Computing, Networking, Storage and Analysis, pp. 588 1-15, 2022. URL https://api.semanticscholar.org/CorpusID:250243681. 589
- Marco Ancona, Enea Ceolini, Cengiz Öztireli, and Markus Gross. Gradient-based attribution methods. *Explainable AI: Interpreting, explaining and visualizing deep learning*, pp. 169–191, 2019. 592
- Anthropic. The claude 3 model family: Opus, sonnet, haiku, 2024. URL https://api. semanticscholar.org/CorpusID:268232499.

- Jasmijn Bastings, Sebastian Ebert, Polina Zablotskaia, Anders Sandholm, and Katja Filippova. "will 595 you find these shortcuts?" a protocol for evaluating the faithfulness of input salience methods for 596 text classification. In Yoav Goldberg, Zornitsa Kozareva, and Yue Zhang (eds.), Proceedings of 597 the 2022 Conference on Empirical Methods in Natural Language Processing, pp. 976–991, Abu 598 Dhabi, United Arab Emirates, December 2022. Association for Computational Linguistics. doi: 10.18653/v1/2022.emnlp-main.64. URL https://aclanthology.org/2022.emnlp-main.64. 600 David Bau. Baukit, 2022. URL https://github.com/davidbau/baukit. 601 602 David Bau, Bolei Zhou, Aditya Khosla, Aude Oliva, and Antonio Torralba. Network dissection: 603 Quantifying interpretability of deep visual representations. In *Proceedings of the IEEE Conference* 604 on Computer Vision and Pattern Recognition (CVPR), pp. 6541–6549, July 2017. 605 Stella Biderman, Hailey Schoelkopf, Quentin Gregory Anthony, Herbie Bradley, Kyle O'Brien, 606 Eric Hallahan, Mohammad Aflah Khan, Shivanshu Purohit, Usvsn Sai Prashanth, Edward Raff, 607 Aviya Skowron, Lintang Sutawika, and Oskar Van Der Wal. Pythia: A suite for analyzing large 608 language models across training and scaling. In Andreas Krause, Emma Brunskill, Kyunghyun 609 Cho, Barbara Engelhardt, Sivan Sabato, and Jonathan Scarlett (eds.), Proceedings of the 40th 610 International Conference on Machine Learning, volume 202 of Proceedings of Machine Learning 611 Research, pp. 2397–2430. PMLR, 23–29 Jul 2023. URL https://proceedings.mlr.press/ 612 v202/biderman23a.html. 613 Alexander Borzunov, Dmitry Baranchuk, Tim Dettmers, Maksim Riabinin, Younes Belkada, Artem 614 Chumachenko, Pavel Samygin, and Colin Raffel. Petals: Collaborative inference and fine-tuning 615 of large models. In Danushka Bollegala, Ruihong Huang, and Alan Ritter (eds.), Proceedings 616 of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 3: System 617 Demonstrations), pp. 558–568, Toronto, Canada, July 2023. Association for Computational 618 Linguistics. doi: 10.18653/v1/2023.acl-demo.54. URL https://aclanthology.org/2023. 619 acl-demo.54. 620 621 Léon Bottou and Patrick Gallinari. A framework for the cooperation of learning algorithms. Advances 622 in neural information processing systems, 3, 1990. 623 Léon Bottou and Yann Le Cun. Sn: A simulator for connectionist models. In Proceedings of 624 NeuroNimes 88, pp. 371–382, Nimes, France, 1988. URL http://leon.bottou.org/papers/ 625 bottou-lecun-88. 626 627 James Bradbury, Roy Frostig, Peter Hawkins, Matthew James Johnson, Chris Leary, Dougal 628 Maclaurin, George Necula, Adam Paszke, Jake VanderPlas, Skye Wanderman-Milne, and 629 Qiao Zhang. JAX: composable transformations of Python+NumPy programs, 2018. URL 630 http://github.com/google/jax. 631 Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, 632 Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language models are 633 few-shot learners. Advances in neural information processing systems, 33:1877–1901, 2020. 634 635 Benjamin S Bucknall and Robert F Trager. Structured access for third-party research on frontier 636 ai models: Investigating researcher's model access requirements, 2023. URL https://cdn. 637 governance.ai/Structured_Access_for_Third-Party_Research.pdf. 638 Stephen Casper, Carson Ezell, Charlotte Siegmann, Noam Kolt, Taylor Lynn Curtis, Benjamin 639 Bucknall, Andreas Haupt, Kevin Wei, Jérémy Scheurer, Marius Hobbhahn, et al. Black-box access 640 is insufficient for rigorous ai audits. In The 2024 ACM Conference on Fairness, Accountability, 641 and Transparency, pp. 2254-2272, 2024. 642 643
- Lequn Chen, Zihao Ye, Yongji Wu, Danyang Zhuo, Luis Ceze, and Arvind Krishnamurthy. Punica: Multi-tenant lora serving. *Proceedings of Machine Learning and Systems*, 6:1–13, 2024.
- Olivier Dehaene, Nicolas Patry, David Richard Holtz, Daniël de Kok, and contribu tors. HuggingFace text generation interface. URL https://github.com/huggingface/
 text-generation-inference, 2024.

Tim Dettmers, Artidoro Pagnoni, Ari Holtzman, and Luke Zettlemoyer. QLoRA: Efficient finetuning of quantized LLMs. In *Thirty-seventh Conference on Neural Information Processing Systems*, 2023. URL https://openreview.net/forum?id=OUIFPHEgJU.

650 651

648

649

652 Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha 653 Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, Anirudh Goyal, Anthony Hartshorn, 654 Aobo Yang, Archi Mitra, Archie Sravankumar, Artem Korenev, Arthur Hinsvark, Arun Rao, Aston Zhang, Aurelien Rodriguez, Austen Gregerson, Ava Spataru, Baptiste Roziere, Bethany Biron, 655 Binh Tang, Bobbie Chern, Charlotte Caucheteux, Chaya Nayak, Chloe Bi, Chris Marra, Chris 656 McConnell, Christian Keller, Christophe Touret, Chunyang Wu, Corinne Wong, Cristian Canton 657 Ferrer, Cyrus Nikolaidis, Damien Allonsius, Daniel Song, Danielle Pintz, Danny Livshits, David 658 Esiobu, Dhruv Choudhary, Dhruv Mahajan, Diego Garcia-Olano, Diego Perino, Dieuwke Hupkes, 659 Egor Lakomkin, Ehab AlBadawy, Elina Lobanova, Emily Dinan, Eric Michael Smith, Filip Radenovic, Frank Zhang, Gabriel Synnaeve, Gabrielle Lee, Georgia Lewis Anderson, Graeme 661 Nail, Gregoire Mialon, Guan Pang, Guillem Cucurell, Hailey Nguyen, Hannah Korevaar, Hu Xu, 662 Hugo Touvron, Iliyan Zarov, Imanol Arrieta Ibarra, Isabel Kloumann, Ishan Misra, Ivan Evtimov, 663 Jade Copet, Jaewon Lee, Jan Geffert, Jana Vranes, Jason Park, Jay Mahadeokar, Jeet Shah, Jelmer van der Linde, Jennifer Billock, Jenny Hong, Jenya Lee, Jeremy Fu, Jianfeng Chi, Jianyu 665 Huang, Jiawen Liu, Jie Wang, Jiecao Yu, Joanna Bitton, Joe Spisak, Jongsoo Park, Joseph Rocca, Joshua Johnstun, Joshua Saxe, Junteng Jia, Kalyan Vasuden Alwala, Kartikeya Upasani, Kate Plawiak, Ke Li, Kenneth Heafield, Kevin Stone, Khalid El-Arini, Krithika Iyer, Kshitiz 667 Malik, Kuenley Chiu, Kunal Bhalla, Lauren Rantala-Yeary, Laurens van der Maaten, Lawrence 668 Chen, Liang Tan, Liz Jenkins, Louis Martin, Lovish Madaan, Lubo Malo, Lukas Blecher, Lukas 669 Landzaat, Luke de Oliveira, Madeline Muzzi, Mahesh Pasupuleti, Mannat Singh, Manohar Paluri, 670 Marcin Kardas, Mathew Oldham, Mathieu Rita, Maya Pavlova, Melanie Kambadur, Mike Lewis, 671 Min Si, Mitesh Kumar Singh, Mona Hassan, Naman Goyal, Narjes Torabi, Nikolay Bashlykov, 672 Nikolay Bogoychev, Niladri Chatterji, Olivier Duchenne, Onur Çelebi, Patrick Alrassy, Pengchuan 673 Zhang, Pengwei Li, Petar Vasic, Peter Weng, Prajjwal Bhargava, Pratik Dubal, Praveen Krishnan, 674 Punit Singh Koura, Puxin Xu, Qing He, Qingxiao Dong, Ragavan Srinivasan, Raj Ganapathy, 675 Ramon Calderer, Ricardo Silveira Cabral, Robert Stojnic, Roberta Raileanu, Rohit Girdhar, Rohit 676 Patel, Romain Sauvestre, Ronnie Polidoro, Roshan Sumbaly, Ross Taylor, Ruan Silva, Rui Hou, Rui Wang, Saghar Hosseini, Sahana Chennabasappa, Sanjay Singh, Sean Bell, Seohyun Sonia 677 Kim, Sergey Edunov, Shaoliang Nie, Sharan Narang, Sharath Raparthy, Sheng Shen, Shengye Wan, 678 Shruti Bhosale, Shun Zhang, Simon Vandenhende, Soumya Batra, Spencer Whitman, Sten Sootla, 679 Stephane Collot, Suchin Gururangan, Sydney Borodinsky, Tamar Herman, Tara Fowler, Tarek 680 Sheasha, Thomas Georgiou, Thomas Scialom, Tobias Speckbacher, Todor Mihaylov, Tong Xiao, Ujjwal Karn, Vedanuj Goswami, Vibhor Gupta, Vignesh Ramanathan, Viktor Kerkez, Vincent 682 Gonguet, Virginie Do, Vish Vogeti, Vladan Petrovic, Weiwei Chu, Wenhan Xiong, Wenyin Fu, Whitney Meers, Xavier Martinet, Xiaodong Wang, Xiaoqing Ellen Tan, Xinfeng Xie, Xuchao Jia, 684 Xuewei Wang, Yaelle Goldschlag, Yashesh Gaur, Yasmine Babaei, Yi Wen, Yiwen Song, Yuchen 685 Zhang, Yue Li, Yuning Mao, Zacharie Delpierre Coudert, Zheng Yan, Zhengxing Chen, Zoe 686 Papakipos, Aaditya Singh, Aaron Grattafiori, Abha Jain, Adam Kelsey, Adam Shajnfeld, Adithya 687 Gangidi, Adolfo Victoria, Ahuva Goldstand, Ajay Menon, Ajay Sharma, Alex Boesenberg, Alex Vaughan, Alexei Baevski, Allie Feinstein, Amanda Kallet, Amit Sangani, Anam Yunus, Andrei 688 Lupu, Andres Alvarado, Andrew Caples, Andrew Gu, Andrew Ho, Andrew Poulton, Andrew 689 Ryan, Ankit Ramchandani, Annie Franco, Aparajita Saraf, Arkabandhu Chowdhury, Ashley 690 Gabriel, Ashwin Bharambe, Assaf Eisenman, Azadeh Yazdan, Beau James, Ben Maurer, Benjamin 691 Leonhardi, Bernie Huang, Beth Loyd, Beto De Paola, Bhargavi Paranjape, Bing Liu, Bo Wu, 692 Boyu Ni, Braden Hancock, Bram Wasti, Brandon Spence, Brani Stojkovic, Brian Gamido, Britt 693 Montalvo, Carl Parker, Carly Burton, Catalina Mejia, Changhan Wang, Changkyu Kim, Chao Zhou, Chester Hu, Ching-Hsiang Chu, Chris Cai, Chris Tindal, Christoph Feichtenhofer, Damon Civin, Dana Beaty, Daniel Kreymer, Daniel Li, Danny Wyatt, David Adkins, David Xu, Davide Testuggine, Delia David, Devi Parikh, Diana Liskovich, Didem Foss, Dingkang Wang, Duc Le, 697 Dustin Holland, Edward Dowling, Eissa Jamil, Elaine Montgomery, Eleonora Presani, Emily Hahn, Emily Wood, Erik Brinkman, Esteban Arcaute, Evan Dunbar, Evan Smothers, Fei Sun, Felix Kreuk, Feng Tian, Firat Ozgenel, Francesco Caggioni, Francisco Guzmán, Frank Kanayet, Frank 699 Seide, Gabriela Medina Florez, Gabriella Schwarz, Gada Badeer, Georgia Swee, Gil Halpern, Govind Thattai, Grant Herman, Grigory Sizov, Guangyi, Zhang, Guna Lakshminarayanan, Hamid Shojanazeri, Han Zou, Hannah Wang, Hanwen Zha, Haroun Habeeb, Harrison Rudolph, Helen

702 Suk, Henry Aspegren, Hunter Goldman, Ibrahim Damlaj, Igor Molybog, Igor Tufanov, Irina-703 Elena Veliche, Itai Gat, Jake Weissman, James Geboski, James Kohli, Japhet Asher, Jean-Baptiste 704 Gaya, Jeff Marcus, Jeff Tang, Jennifer Chan, Jenny Zhen, Jeremy Reizenstein, Jeremy Teboul, 705 Jessica Zhong, Jian Jin, Jingyi Yang, Joe Cummings, Jon Carvill, Jon Shepard, Jonathan McPhie, 706 Jonathan Torres, Josh Ginsburg, Junjie Wang, Kai Wu, Kam Hou U, Karan Saxena, Karthik Prasad, Kartikay Khandelwal, Katayoun Zand, Kathy Matosich, Kaushik Veeraraghavan, Kelly Michelena, Keqian Li, Kun Huang, Kunal Chawla, Kushal Lakhotia, Kyle Huang, Lailin Chen, 708 Lakshya Garg, Lavender A, Leandro Silva, Lee Bell, Lei Zhang, Liangpeng Guo, Licheng Yu, 709 Liron Moshkovich, Luca Wehrstedt, Madian Khabsa, Manav Avalani, Manish Bhatt, Maria 710 Tsimpoukelli, Martynas Mankus, Matan Hasson, Matthew Lennie, Matthias Reso, Maxim Groshev, 711 Maxim Naumov, Maya Lathi, Meghan Keneally, Michael L. Seltzer, Michal Valko, Michelle 712 Restrepo, Mihir Patel, Mik Vyatskov, Mikayel Samvelyan, Mike Clark, Mike Macey, Mike Wang, 713 Miquel Jubert Hermoso, Mo Metanat, Mohammad Rastegari, Munish Bansal, Nandhini Santhanam, 714 Natascha Parks, Natasha White, Navyata Bawa, Nayan Singhal, Nick Egebo, Nicolas Usunier, 715 Nikolay Pavlovich Laptev, Ning Dong, Ning Zhang, Norman Cheng, Oleg Chernoguz, Olivia 716 Hart, Omkar Salpekar, Ozlem Kalinli, Parkin Kent, Parth Parekh, Paul Saab, Pavan Balaji, Pedro 717 Rittner, Philip Bontrager, Pierre Roux, Piotr Dollar, Polina Zvyagina, Prashant Ratanchandani, Pritish Yuvraj, Qian Liang, Rachad Alao, Rachel Rodriguez, Rafi Ayub, Raghotham Murthy, 718 Raghu Nayani, Rahul Mitra, Raymond Li, Rebekkah Hogan, Robin Battey, Rocky Wang, Rohan 719 Maheswari, Russ Howes, Ruty Rinott, Sai Jayesh Bondu, Samyak Datta, Sara Chugh, Sara 720 Hunt, Sargun Dhillon, Sasha Sidorov, Satadru Pan, Saurabh Verma, Seiji Yamamoto, Sharadh 721 Ramaswamy, Shaun Lindsay, Shaun Lindsay, Sheng Feng, Shenghao Lin, Shengxin Cindy Zha, 722 Shiva Shankar, Shuqiang Zhang, Shuqiang Zhang, Sinong Wang, Sneha Agarwal, Soji Sajuyigbe, 723 Soumith Chintala, Stephanie Max, Stephen Chen, Steve Kehoe, Steve Satterfield, Sudarshan 724 Govindaprasad, Sumit Gupta, Sungmin Cho, Sunny Virk, Suraj Subramanian, Sy Choudhury, 725 Sydney Goldman, Tal Remez, Tamar Glaser, Tamara Best, Thilo Kohler, Thomas Robinson, Tianhe 726 Li, Tianjun Zhang, Tim Matthews, Timothy Chou, Tzook Shaked, Varun Vontimitta, Victoria Ajayi, 727 Victoria Montanez, Vijai Mohan, Vinay Satish Kumar, Vishal Mangla, Vítor Albiero, Vlad Ionescu, Vlad Poenaru, Vlad Tiberiu Mihailescu, Vladimir Ivanov, Wei Li, Wenchen Wang, Wenwen Jiang, 728 Wes Bouaziz, Will Constable, Xiaocheng Tang, Xiaofang Wang, Xiaojian Wu, Xiaolan Wang, 729 Xide Xia, Xilun Wu, Xinbo Gao, Yanjun Chen, Ye Hu, Ye Jia, Ye Qi, Yenda Li, Yilin Zhang, 730 Ying Zhang, Yossi Adi, Youngjin Nam, Yu, Wang, Yuchen Hao, Yundi Qian, Yuzi He, Zach Rait, 731 Zachary DeVito, Zef Rosnbrick, Zhaoduo Wen, Zhenyu Yang, and Zhiwei Zhao. The llama 3 herd 732 of models, 2024. URL https://arxiv.org/abs/2407.21783. 733

- Nelson Elhage, Neel Nanda, Catherine Olsson, Tom Henighan, Nicholas Joseph, Ben Mann, Amanda Askell, Yuntao Bai, Anna Chen, Tom Conerly, Nova DasSarma, Dawn Drain, Deep Ganguli, Zac Hatfield-Dodds, Danny Hernandez, Andy Jones, Jackson Kernion, Liane Lovitt, Kamal Ndousse, Dario Amodei, Tom Brown, Jack Clark, Jared Kaplan, Sam McCandlish, and Chris Olah. Garcon. URL https://transformer-circuits.pub/2021/garcon/index.html, 2021.

742

- Javier Ferrando, Gabriele Sarti, Arianna Bisazza, and Marta R Costa-jussà. A primer on the inner workings of transformer-based language models. *arXiv preprint arXiv:2405.00208*, 2024.
- 744 Gemini Team, Machel Reid, Nikolay Savinov, Denis Teplyashin, Dmitry Lepikhin, Timothy P. Lilli-745 crap, Jean-Baptiste Alayrac, Radu Soricut, Angeliki Lazaridou, Orhan Firat, Julian Schrittwieser, 746 Ioannis Antonoglou, Rohan Anil, Sebastian Borgeaud, Andrew M. Dai, Katie Millican, Ethan 747 Dyer, Mia Glaese, Thibault Sottiaux, Benjamin Lee, Fabio Viola, Malcolm Reynolds, Yuanzhong 748 Xu, James Molloy, Jilin Chen, Michael Isard, Paul Barham, Tom Hennigan, Ross McIlroy, Melvin 749 Johnson, Johan Schalkwyk, Eli Collins, Eliza Rutherford, Erica Moreira, Kareem W. Ayoub, 750 Megha Goel, Clemens Meyer, Gregory Thornton, Zhen Yang, Henryk Michalewski, Zaheer Ab-751 bas, Nathan Schucher, Ankesh Anand, Richard Ives, James Keeling, Karel Lenc, Salem Haykal, 752 Siamak Shakeri, Pranav Shyam, Aakanksha Chowdhery, Roman Ring, Stephen Spencer, Eren 753 Sezener, Luke Vilnis, Oscar Chang, Nobuyuki Morioka, George Tucker, Ce Zheng, Oliver Woodman, Nithya Attaluri, Tomás Kociský, Evgenii Eltyshev, Xi Chen, Timothy Chung, Vittorio Selo, 754 Siddhartha Brahma, Petko Georgiev, Ambrose Slone, Zhenkai Zhu, James Lottes, Siyuan Qiao, 755 Ben Caine, Sebastian Riedel, Alex Tomala, Martin Chadwick, J Christopher Love, Peter Choy,

756 Sid Mittal, Neil Houlsby, Yunhao Tang, Matthew Lamm, Libin Bai, Qiao Zhang, Luheng He, Yong Cheng, Peter Humphreys, Yujia Li, Sergey Brin, Albin Cassirer, Ying-Qi Miao, Lukás 758 Zilka, Taylor Tobin, Kelvin Xu, Lev Proleev, Daniel Sohn, Alberto Magni, Lisa Anne Hendricks, 759 Isabel Gao, Santiago Ontan'on, Oskar Bunyan, Nathan Byrd, Abhanshu Sharma, Biao Zhang, 760 Mario Pinto, Rishika Sinha, Harsh Mehta, Dawei Jia, Sergi Caelles, Albert Webson, Alex Morris, Becca Roelofs, Yifan Ding, Robin Strudel, Xuehan Xiong, Marvin Ritter, Mostafa Dehghani, 761 Rahma Chaabouni, Abhijit Karmarkar, Guangda Lai, Fabian Mentzer, Bibo Xu, YaGuang Li, 762 Yujing Zhang, Tom Le Paine, Alex Goldin, Behnam Neyshabur, Kate Baumli, Anselm Levskaya, Michael Laskin, Wenhao Jia, Jack W. Rae, Kefan Xiao, Antoine He, Skye Giordano, Laksh-764 man Yagati, Jean-Baptiste Lespiau, Paul Natsev, Sanjay Ganapathy, Fangyu Liu, Danilo Martins, 765 Nanxin Chen, Yunhan Xu, Megan Barnes, Rhys May, Arpi Vezer, Junhyuk Oh, Ken Franko, 766 Sophie Bridgers, Ruizhe Zhao, Boxi Wu, Basil Mustafa, Sean Sechrist, Emilio Parisotto, Thanu-767 malayan Sankaranarayana Pillai, Chris Larkin, Chenjie Gu, Christina Sorokin, Maxim Krikun, 768 Alexey Guseynov, Jessica Landon, Romina Datta, Alexander Pritzel, Phoebe Thacker, Fan Yang, 769 Kevin Hui, A.E. Hauth, Chih-Kuan Yeh, David Barker, Justin Mao-Jones, Sophia Austin, Hannah 770 Sheahan, Parker Schuh, James Svensson, Rohan Jain, Vinay Venkatesh Ramasesh, Anton Briukhov, Da-Woon Chung, Tamara von Glehn, Christina Butterfield, Priya Jhakra, Matt Wiethoff, Justin 771 Frye, Jordan Grimstad, Beer Changpinyo, Charline Le Lan, Anna Bortsova, Yonghui Wu, Paul 772 Voigtlaender, Tara N. Sainath, Charlotte Smith, Will Hawkins, Kris Cao, James Besley, Srivatsan 773 Srinivasan, Mark Omernick, Colin Gaffney, Gabriela de Castro Surita, Ryan Burnell, Bogdan 774 Damoc, Junwhan Ahn, Andrew Brock, Mantas Pajarskas, Anastasia Petrushkina, Seb Noury, 775 Lorenzo Blanco, Kevin Swersky, Arun Ahuja, Thi Avrahami, Vedant Misra, Raoul de Liedekerke, 776 Mariko Iinuma, Alex Polozov, Sarah York, George van den Driessche, Paul Michel, Justin Chiu, 777 Rory Blevins, Zach Gleicher, Adrià Recasens, Alban Rrustemi, Elena Gribovskaya, Aurko Roy, 778 Wiktor Gworek, S'ebastien M. R. Arnold, Lisa Lee, James Lee-Thorp, Marcello Maggioni, Demis 779 Hassabis, Koray Kavukcuoglu, Jeffrey Dean, and Oriol Vinyals. Gemini 1.5: Unlocking multimodal understanding across millions of tokens of context. ArXiv, abs/2403.05530, 2024. URL 781 https://api.semanticscholar.org/CorpusID:268297180. 782

- Mor Geva, Jasmijn Bastings, Katja Filippova, and Amir Globerson. Dissecting recall of factual 783 associations in auto-regressive language models. In Proceedings of the 2023 Conference on 784 Empirical Methods in Natural Language Processing, pp. 12216–12235, 2023. 785
- 786 Dan Hendrycks, Collin Burns, Steven Basart, Andy Zou, Mantas Mazeika, Dawn Song, and Jacob 787 Steinhardt. Measuring massive multitask language understanding. In International Conference on 788 Learning Representations, 2021.

789

790

791

792

797

- J. Edward Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, and Weizhu Chen. Lora: Low-rank adaptation of large language models. ArXiv, abs/2106.09685, 2021. URL https://api.semanticscholar.org/CorpusID:235458009.
- 793 Robert Huben, Hoagy Cunningham, Logan Riggs Smith, Aidan Ewart, and Lee Sharkey. Sparse 794 autoencoders find highly interpretable features in language models. In The Twelfth International Conference on Learning Representations, 2024. URL https://openreview.net/forum?id= 795 F76bwRSLeK. 796
- Huggingface. Huggingface open llm leaderboard archive, 2024. URL https://huggingface.co/ 798 spaces/open-llm-leaderboard-old/open_llm_leaderboard. Last Accessed: 2024-09-30.
- 800 Albert Qiaochu Jiang, Alexandre Sablayrolles, Arthur Mensch, Chris Bamford, Devendra Singh Chaplot, Diego de Las Casas, Florian Bressand, Gianna Lengyel, Guillaume Lample, Lucile 801 Saulnier, L'elio Renard Lavaud, Marie-Anne Lachaux, Pierre Stock, Teven Le Scao, Thibaut Lavril, 802 Thomas Wang, Timothée Lacroix, and William El Sayed. Mistral 7b. ArXiv, abs/2310.06825, 803 2023. URL https://api.semanticscholar.org/CorpusID:263830494. 804
- 805 Daniel D. Johnson. Penzai + Treescope: A toolkit for interpreting, visualizing, and editing models as 806 data. ICML 2024 Workshop on Mechanistic Interpretability, 2024. URL https://openreview. 807 net/forum?id=KVSgEXrMDH. 808
- János Kramár, Tom Lieberum, Rohin Shah, and Neel Nanda. Atp*: An efficient and scalable method for localizing llm behaviour to components. arXiv preprint arXiv:2403.00745, 2024.

827

- Woosuk Kwon, Zhuohan Li, Siyuan Zhuang, Ying Sheng, Lianmin Zheng, Cody Hao Yu, Joseph E.
 Gonzalez, Hao Zhang, and Ion Stoica. Efficient memory management for large language model
 serving with pagedattention. In *Proceedings of the ACM SIGOPS 29th Symposium on Operating Systems Principles*, 2023.
- Brian Lester, Rami Al-Rfou, and Noah Constant. The power of scale for parameter-efficient prompt tuning. In Marie-Francine Moens, Xuanjing Huang, Lucia Specia, and Scott Wentau Yih (eds.), *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pp. 3045–3059, Online and Punta Cana, Dominican Republic, November 2021. Association for Computational Linguistics. doi: 10.18653/v1/2021.emnlp-main.243. URL https://aclanthology.org/2021.emnlp-main.243.
- Kenneth Li, Oam Patel, Fernanda Viégas, Hanspeter Pfister, and Martin Wattenberg. Inference-time
 intervention: Eliciting truthful answers from a language model. *Advances in Neural Information Processing Systems*, 36, 2024.
- Samuel Marks, Can Rager, Eric J. Michaud, Yonatan Belinkov, David Bau, and Aaron Mueller.
 Sparse feature circuits: Discovering and editing interpretable causal graphs in language models, 2024. URL https://arxiv.org/abs/2403.19647.
- Kevin Meng, David Bau, Alex Andonian, and Yonatan Belinkov. Locating and editing factual associations in gpt. In *Neural Information Processing Systems*, 2022. URL https://openreview.net/forum?id=-h6WAS6eE4.
- Neel Nanda and Joseph Bloom. Transformerlens. https://github.com/TransformerLensOrg/
 TransformerLens, 2022.
- 833
 834 Nvidia. TensorRT-LLM. URL https://github.com/NVIDIA/TensorRT-LLM, 2024.

835 OpenAI, Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni 836 Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, Red Avila, Igor 837 Babuschkin, Suchir Balaji, Valerie Balcom, Paul Baltescu, Haiming Bao, Mo Bavarian, Jeff 838 Belgum, Irwan Bello, Jake Berdine, Gabriel Bernadett-Shapiro, Christopher Berner, Lenny Bogdonoff, Oleg Boiko, Madelaine Boyd, Anna-Luisa Brakman, Greg Brockman, Tim Brooks, Miles 839 Brundage, Kevin Button, Trevor Cai, Rosie Campbell, Andrew Cann, Brittany Carey, Chelsea 840 Carlson, Rory Carmichael, Brooke Chan, Che Chang, Fotis Chantzis, Derek Chen, Sully Chen, 841 Ruby Chen, Jason Chen, Mark Chen, Benjamin Chess, Chester Cho, Casey Chu, Hyung Won 842 Chung, Dave Cummings, Jeremiah Currier, Yunxing Dai, Cory Decareaux, Thomas Degry, Noah 843 Deutsch, Damien Deville, Arka Dhar, David Dohan, Steve Dowling, Sheila Dunning, Adrien 844 Ecoffet, Atty Eleti, Tyna Eloundou, David Farhi, Liam Fedus, Niko Felix, Sim'on Posada Fishman, 845 Juston Forte, Isabella Fulford, Leo Gao, Elie Georges, Christian Gibson, Vik Goel, Tarun Gogineni, 846 Gabriel Goh, Raphael Gontijo-Lopes, Jonathan Gordon, Morgan Grafstein, Scott Gray, Ryan 847 Greene, Joshua Gross, Shixiang Shane Gu, Yufei Guo, Chris Hallacy, Jesse Han, Jeff Harris, 848 Yuchen He, Mike Heaton, Johannes Heidecke, Chris Hesse, Alan Hickey, Wade Hickey, Peter 849 Hoeschele, Brandon Houghton, Kenny Hsu, Shengli Hu, Xin Hu, Joost Huizinga, Shantanu Jain, Shawn Jain, Joanne Jang, Angela Jiang, Roger Jiang, Haozhun Jin, Denny Jin, Shino Jomoto, Billie 850 Jonn, Heewoo Jun, Tomer Kaftan, Lukasz Kaiser, Ali Kamali, Ingmar Kanitscheider, Nitish Shirish 851 Keskar, Tabarak Khan, Logan Kilpatrick, Jong Wook Kim, Christina Kim, Yongjik Kim, Hendrik 852 Kirchner, Jamie Ryan Kiros, Matthew Knight, Daniel Kokotajlo, Lukasz Kondraciuk, Andrew 853 Kondrich, Aris Konstantinidis, Kyle Kosic, Gretchen Krueger, Vishal Kuo, Michael Lampe, Ikai 854 Lan, Teddy Lee, Jan Leike, Jade Leung, Daniel Levy, Chak Ming Li, Rachel Lim, Molly Lin, 855 Stephanie Lin, Mateusz Litwin, Theresa Lopez, Ryan Lowe, Patricia Lue, Anna Adeola Makanju, Kim Malfacini, Sam Manning, Todor Markov, Yaniv Markovski, Bianca Martin, Katie Mayer, Andrew Mayne, Bob McGrew, Scott Mayer McKinney, Christine McLeavey, Paul McMillan, 858 Jake McNeil, David Medina, Aalok Mehta, Jacob Menick, Luke Metz, Andrey Mishchenko, 859 Pamela Mishkin, Vinnie Monaco, Evan Morikawa, Daniel P. Mossing, Tong Mu, Mira Murati, Oleg Murk, David M'ely, Ashvin Nair, Reiichiro Nakano, Rajeev Nayak, Arvind Neelakantan, Richard Ngo, Hyeonwoo Noh, Ouyang Long, Cullen O'Keefe, Jakub W. Pachocki, Alex Paino, 861 Joe Palermo, Ashley Pantuliano, Giambattista Parascandolo, Joel Parish, Emy Parparita, Alexandre 862 Passos, Mikhail Pavlov, Andrew Peng, Adam Perelman, Filipe de Avila Belbute Peres, Michael Petrov, Henrique Pondé de Oliveira Pinto, Michael Pokorny, Michelle Pokrass, Vitchyr H. Pong,

864 Tolly Powell, Alethea Power, Boris Power, Elizabeth Proehl, Raul Puri, Alec Radford, Jack 865 Rae, Aditya Ramesh, Cameron Raymond, Francis Real, Kendra Rimbach, Carl Ross, Bob Rot-866 sted, Henri Roussez, Nick Ryder, Mario D. Saltarelli, Ted Sanders, Shibani Santurkar, Girish 867 Sastry, Heather Schmidt, David Schnurr, John Schulman, Daniel Selsam, Kyla Sheppard, Toki 868 Sherbakov, Jessica Shieh, Sarah Shoker, Pranav Shyam, Szymon Sidor, Eric Sigler, Maddie Simens, Jordan Sitkin, Katarina Slama, Ian Sohl, Benjamin D. Sokolowsky, Yang Song, Natalie Staudacher, Felipe Petroski Such, Natalie Summers, Ilya Sutskever, Jie Tang, Nikolas A. Tezak, 870 Madeleine Thompson, Phil Tillet, Amin Tootoonchian, Elizabeth Tseng, Preston Tuggle, Nick 871 Turley, Jerry Tworek, Juan Felipe Cer'on Uribe, Andrea Vallone, Arun Vijayvergiya, Chelsea 872 Voss, Carroll L. Wainwright, Justin Jay Wang, Alvin Wang, Ben Wang, Jonathan Ward, Jason 873 Wei, CJ Weinmann, Akila Welihinda, Peter Welinder, Jiayi Weng, Lilian Weng, Matt Wiethoff, 874 Dave Willner, Clemens Winter, Samuel Wolrich, Hannah Wong, Lauren Workman, Sherwin Wu, 875 Jeff Wu, Michael Wu, Kai Xiao, Tao Xu, Sarah Yoo, Kevin Yu, Qiming Yuan, Wojciech Zaremba, 876 Rowan Zellers, Chong Zhang, Marvin Zhang, Shengjia Zhao, Tianhao Zheng, Juntang Zhuang, 877 William Zhuk, and Barret Zoph. Gpt-4 technical report. arXiv preprint arXiv:2303.08774, 2023. 878 URL https://api.semanticscholar.org/CorpusID:257532815.

- PaperswithCode. Multi-task language understanding on mmlu, 2024. URL https: //paperswithcode.com/sota/multi-task-language-understanding-on-mmlu. Last Accessed: 2024-09-30.
- Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor
 Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, et al. Pytorch: An imperative style,
 high-performance deep learning library. *Advances in neural information processing systems*, 32, 2019.
- Roma Patel and Ellie Pavlick. Mapping language models to grounded conceptual spaces. In International Conference on Learning Representations, 2022. URL https://openreview.net/ forum?id=gJcEM8sxHK.
- Alec Radford, Jeff Wu, Rewon Child, David Luan, Dario Amodei, and Ilya Sutskever. Language
 models are unsupervised multitask learners. *OpenAI Blog*, 2019.
- 893 Jack W. Rae, Sebastian Borgeaud, Trevor Cai, Katie Millican, Jordan Hoffmann, Francis Song, 894 John Aslanides, Sarah Henderson, Roman Ring, Susannah Young, Eliza Rutherford, Tom Hen-895 nigan, Jacob Menick, Albin Cassirer, Richard Powell, George van den Driessche, Lisa Anne Hendricks, Maribeth Rauh, Po-Sen Huang, Amelia Glaese, Johannes Welbl, Sumanth Dathathri, 896 Saffron Huang, Jonathan Uesato, John F. J. Mellor, Irina Higgins, Antonia Creswell, Nathan 897 McAleese, Amy Wu, Erich Elsen, Siddhant M. Jayakumar, Elena Buchatskaya, David Budden, Esme Sutherland, Karen Simonyan, Michela Paganini, L. Sifre, Lena Martens, Xiang Lor-899 raine Li, Adhiguna Kuncoro, Aida Nematzadeh, Elena Gribovskaya, Domenic Donato, Angeliki 900 Lazaridou, Arthur Mensch, Jean-Baptiste Lespiau, Maria Tsimpoukelli, N. K. Grigorev, Doug 901 Fritz, Thibault Sottiaux, Mantas Pajarskas, Tobias Pohlen, Zhitao Gong, Daniel Toyama, Cy-902 prien de Masson d'Autume, Yujia Li, Tayfun Terzi, Vladimir Mikulik, Igor Babuschkin, Aidan 903 Clark, Diego de Las Casas, Aurelia Guy, Chris Jones, James Bradbury, Matthew G. Johnson, 904 Blake A. Hechtman, Laura Weidinger, Iason Gabriel, William S. Isaac, Edward Lockhart, Si-905 mon Osindero, Laura Rimell, Chris Dyer, Oriol Vinyals, Kareem W. Ayoub, Jeff Stanway, L. L. 906 Bennett, Demis Hassabis, Koray Kavukcuoglu, and Geoffrey Irving. Scaling language models: Methods, analysis & insights from training gopher. ArXiv, abs/2112.11446, 2021. URL 907 https://api.semanticscholar.org/CorpusID:245353475. 908
- 909 Teven Le Scao, Angela Fan, Christopher Akiki, Ellie Pavlick, Suzana Ili'c, Daniel Hesslow, Roman 910 Castagn'e, Alexandra Sasha Luccioni, François Yvon, Matthias Gallé, Jonathan Tow, Alexander M. 911 Rush, Stella Biderman, Albert Webson, Pawan Sasanka Ammanamanchi, Thomas Wang, Benoît 912 Sagot, Niklas Muennighoff, Albert Villanova del Moral, Olatunji Ruwase, Rachel Bawden, Stas 913 Bekman, Angelina McMillan-Major, Iz Beltagy, Huu Nguyen, Lucile Saulnier, Samson Tan, Pe-914 dro Ortiz Suarez, Victor Sanh, Hugo Laurenccon, Yacine Jernite, Julien Launay, Margaret Mitchell, 915 Colin Raffel, Aaron Gokaslan, Adi Simhi, Aitor Soroa Etxabe, Alham Fikri Aji, Amit Alfassy, Anna Rogers, Ariel Kreisberg Nitzav, Canwen Xu, Chenghao Mou, Chris C. Emezue, Christo-916 pher Klamm, Colin Leong, Daniel Alexander van Strien, David Ifeoluwa Adelani, Dragomir R. 917 Radev, Eduardo Gonz'alez Ponferrada, Efrat Levkovizh, Ethan Kim, Eyal Natan, Francesco De

918 Toni, Gérard Dupont, Germán Kruszewski, Giada Pistilli, Hady ElSahar, Hamza Benyamina, 919 Hieu Trung Tran, Ian Yu, Idris Abdulmumin, Isaac Johnson, Itziar Gonzalez-Dios, Javier de la 920 Rosa, Jenny Chim, Jesse Dodge, Jian Zhu, Jonathan Chang, Jorg Frohberg, Josephine Tobing, 921 Joydeep Bhattacharjee, Khalid Almubarak, Kimbo Chen, Kyle Lo, Leandro von Werra, Leon 922 Weber, Long Phan, Loubna Ben Allal, Ludovic Tanguy, Manan Dey, Manuel Romero Muñoz, Maraim Masoud, María Grandury, Mario vSavsko, Max Huang, Maximin Coavoux, Mayank 923 Singh, Mike Tian-Jian Jiang, Minh Chien Vu, Mohammad A. Jauhar, Mustafa Ghaleb, Nishant 924 Subramani, Nora Kassner, Nurulaqilla Khamis, Olivier Nguyen, Omar Espejel, Ona de Gibert, 925 Paulo Villegas, Peter Henderson, Pierre Colombo, Priscilla Amuok, Quentin Lhoest, Rheza Harli-926 man, Rishi Bommasani, Roberto L'opez, Rui Ribeiro, Salomey Osei, Sampo Pyysalo, Sebastian 927 Nagel, Shamik Bose, Shamsuddeen Hassan Muhammad, Shanya Sharma, S. Longpre, Somaieh 928 Nikpoor, S. Silberberg, Suhas Pai, Sydney Zink, Tiago Timponi Torrent, Timo Schick, Tristan 929 Thrush, Valentin Danchev, Vassilina Nikoulina, Veronika Laippala, Violette Lepercq, Vrinda 930 Prabhu, Zaid Alyafeai, Zeerak Talat, Arun Raja, Benjamin Heinzerling, Chenglei Si, Elizabeth 931 Salesky, Sabrina J. Mielke, Wilson Y. Lee, Abheesht Sharma, Andrea Santilli, Antoine Chaf-932 fin, Arnaud Stiegler, Debajyoti Datta, Eliza Szczechla, Gunjan Chhablani, Han Wang, Harshit Pandey, Hendrik Strobelt, Jason Alan Fries, Jos Rozen, Leo Gao, Lintang Sutawika, M Saiful Bari, Maged S. Al-Shaibani, Matteo Manica, Nihal V. Nayak, Ryan Teehan, Samuel Albanie, Sheng 934 Shen, Srulik Ben-David, Stephen H. Bach, Taewoon Kim, Tali Bers, Thibault Févry, Trishala 935 Neeraj, Urmish Thakker, Vikas Raunak, Xiang Tang, Zheng-Xin Yong, Zhiqing Sun, Shaked 936 Brody, Y Uri, Hadar Tojarieh, Adam Roberts, Hyung Won Chung, Jaesung Tae, Jason Phang, 937 Ofir Press, Conglong Li, Deepak Narayanan, Hatim Bourfoune, Jared Casper, Jeff Rasley, Max 938 Ryabinin, Mayank Mishra, Minjia Zhang, Mohammad Shoeybi, Myriam Peyrounette, Nicolas 939 Patry, Nouamane Tazi, Omar Sanseviero, Patrick von Platen, Pierre Cornette, Pierre Franccois 940 Lavall'ee, Rémi Lacroix, Samyam Rajbhandari, Sanchit Gandhi, Shaden Smith, Stéphane Requena, 941 Suraj Patil, Tim Dettmers, Ahmed Baruwa, Amanpreet Singh, Anastasia Cheveleva, Anne-Laure 942 Ligozat, Arjun Subramonian, Aur'elie N'ev'eol, Charles Lovering, Daniel H Garrette, Deepak R. 943 Tunuguntla, Ehud Reiter, Ekaterina Taktasheva, Ekaterina Voloshina, Eli Bogdanov, Genta Indra 944 Winata, Hailey Schoelkopf, Jan-Christoph Kalo, Jekaterina Novikova, Jessica Zosa Forde, Xiangru Tang, Jungo Kasai, Ken Kawamura, Liam Hazan, Marine Carpuat, Miruna Clinciu, Najoung Kim, 945 Newton Cheng, Oleg Serikov, Omer Antverg, Oskar van der Wal, Rui Zhang, Ruochen Zhang, 946 Sebastian Gehrmann, Shachar Mirkin, S. Osher Pais, Tatiana Shavrina, Thomas Scialom, Tian Yun, 947 Tomasz Limisiewicz, Verena Rieser, Vitaly Protasov, Vladislav Mikhailov, Yada Pruksachatkun, 948 Yonatan Belinkov, Zachary Bamberger, Zdenvek Kasner, Zdeněk Kasner, Amanda Pestana, Amir 949 Feizpour, Ammar Khan, Amy Faranak, Ananda Santa Rosa Santos, Anthony Hevia, Antigona 950 Unldreaj, Arash Aghagol, Arezoo Abdollahi, Aycha Tammour, Azadeh HajiHosseini, Bahareh 951 Behroozi, Benjamin Ayoade Ajibade, Bharat Kumar Saxena, Carlos Muñoz Ferrandis, Danish 952 Contractor, David M. Lansky, Davis David, Douwe Kiela, Duong Anh Nguyen, Edward Tan, Emi 953 Baylor, Ezinwanne Ozoani, Fatim Tahirah Mirza, Frankline Ononiwu, Habib Rezanejad, H.A. 954 Jones, Indrani Bhattacharya, Irene Solaiman, Irina Sedenko, Isar Nejadgholi, Jan Passmore, Joshua 955 Seltzer, Julio Bonis Sanz, Karen Fort, Lívia Dutra, Mairon Samagaio, Maraim Elbadri, Margot Mieskes, Marissa Gerchick, Martha Akinlolu, Michael McKenna, Mike Qiu, Muhammed Ghauri, 956 Mykola Burynok, Nafis Abrar, Nazneen Rajani, Nour Elkott, Nourhan Fahmy, Olanrewaju Samuel, 957 Ran An, R. P. Kromann, Ryan Hao, Samira Alizadeh, Sarmad Shubber, Silas L. Wang, Sourav 958 Roy, Sylvain Viguier, Thanh-Cong Le, Tobi Oyebade, Trieu Nguyen Hai Le, Yoyo Yang, Zach 959 Nguyen, Abhinav Ramesh Kashyap, Alfredo Palasciano, Alison Callahan, Anima Shukla, Antonio 960 Miranda-Escalada, Ayush Kumar Singh, Benjamin Beilharz, Bo Wang, Caio Matheus Fonseca 961 de Brito, Chenxi Zhou, Chirag Jain, Chuxin Xu, Clémentine Fourrier, Daniel Le'on Perin'an, 962 Daniel Molano, Dian Yu, Enrique Manjavacas, Fabio Barth, Florian Fuhrimann, Gabriel Altay, 963 Giyaseddin Bayrak, Gully Burns, Helena U. Vrabec, Iman I.B. Bello, Isha Dash, Ji Soo Kang, John 964 Giorgi, Jonas Golde, Jose David Posada, Karthi Sivaraman, Lokesh Bulchandani, Lu Liu, Luisa 965 Shinzato, Madeleine Hahn de Bykhovetz, Maiko Takeuchi, Marc Pàmies, María Andrea Castillo, Marianna Nezhurina, Mario Sanger, Matthias Samwald, Michael Cullan, Michael Weinberg, 966 967 M Wolf, Mina Mihaljcic, Minna Liu, Moritz Freidank, Myungsun Kang, Natasha Seelam, Nathan Dahlberg, Nicholas Michio Broad, Nikolaus Muellner, Pascale Fung, Patricia Haller, R. Chan-968 drasekhar, Renata Eisenberg, Robert Martin, Rodrigo Canalli, Rosaline Su, Ruisi Su, Samuel 969 Cahyawijaya, Samuele Garda, Shlok S Deshmukh, Shubhanshu Mishra, Sid Kiblawi, Simon Ott, 970 Sinee Sang-aroonsiri, Srishti Kumar, Stefan Schweter, Sushil Pratap Bharati, Tanmay Laud, Théo Gigant, Tomoya Kainuma, Wojciech Kusa, Yanis Labrak, Yashasvi Bajaj, Y. Venkatraman, Yifan

975

990

991

996

972 Xu, Ying Xu, Yu Xu, Zhee Xao Tan, Zhongli Xie, Zifan Ye, Mathilde Bras, Younes Belkada, 973 and Thomas Wolf. Bloom: A 176b-parameter open-access multilingual language model. ArXiv, abs/2211.05100, 2022. URL https://api.semanticscholar.org/CorpusID:253420279.

- Rylan Schaeffer, Brando Miranda, and Sanmi Koyejo. Are emergent abilities of large language 976 models a mirage? In Thirty-seventh Conference on Neural Information Processing Systems, 2023. 977 URL https://openreview.net/forum?id=ITw9edRDlD. 978
- 979 Ramprasaath R Selvaraju, Michael Cogswell, Abhishek Das, Ramakrishna Vedantam, Devi Parikh, 980 and Dhruv Batra. Grad-cam: Visual explanations from deep networks via gradient-based local-981 ization. In Proceedings of the IEEE international conference on computer vision, pp. 618–626, 982 2017.
- 983 Ying Sheng, Shiyi Cao, Dacheng Li, Coleman Hooper, Nicholas Lee, Shuo Yang, Christo-984 pher Chou, Banghua Zhu, Lianmin Zheng, Kurt Keutzer, Joseph Gonzalez, and Ion Sto-985 ica. Slora: Scalable serving of thousands of lora adapters. In P. Gibbons, G. Pekhi-986 menko, and C. De Sa (eds.), Proceedings of Machine Learning and Systems, volume 6, pp. 987 296-311, 2024. URL https://proceedings.mlsys.org/paper_files/paper/2024/file/ 906419cd502575b617cc489a1a696a67-Paper-Conference.pdf. 989
 - Toby Shevlane. Structured access: an emerging paradigm for safe ai deployment, 2022. URL https://arxiv.org/abs/2201.05159.
- 992 Mukund Sundararajan, Ankur Taly, and Qiqi Yan. Axiomatic attribution for deep networks. In Doina 993 Precup and Yee Whye Teh (eds.), Proceedings of the 34th International Conference on Machine 994 Learning, volume 70 of Proceedings of Machine Learning Research, pp. 3319–3328. PMLR, 995 06-11 Aug 2017. URL https://proceedings.mlr.press/v70/sundararajan17a.html.
- Gemma Team, Thomas Mesnard, Cassidy Hardin, Robert Dadashi, Surya Bhupatiraju, Shreya Pathak, 997 Laurent Sifre, Morgane Rivière, Mihir Sanjay Kale, Juliette Love, Pouya Tafti, Léonard Hussenot, 998 Pier Giuseppe Sessa, Aakanksha Chowdhery, Adam Roberts, Aditya Barua, Alex Botev, Alex 999 Castro-Ros, Ambrose Slone, Amélie Héliou, Andrea Tacchetti, Anna Bulanova, Antonia Paterson, 1000 Beth Tsai, Bobak Shahriari, Charline Le Lan, Christopher A. Choquette-Choo, Clément Crepy, 1001 Daniel Cer, Daphne Ippolito, David Reid, Elena Buchatskaya, Eric Ni, Eric Noland, Geng Yan, 1002 George Tucker, George-Christian Muraru, Grigory Rozhdestvenskiy, Henryk Michalewski, Ian 1003 Tenney, Ivan Grishchenko, Jacob Austin, James Keeling, Jane Labanowski, Jean-Baptiste Lespiau, 1004 Jeff Stanway, Jenny Brennan, Jeremy Chen, Johan Ferret, Justin Chiu, Justin Mao-Jones, Katherine 1005 Lee, Kathy Yu, Katie Millican, Lars Lowe Sjoesund, Lisa Lee, Lucas Dixon, Machel Reid, Maciej Mikuła, Mateo Wirth, Michael Sharman, Nikolai Chinaev, Nithum Thain, Olivier Bachem, Oscar Chang, Oscar Wahltinez, Paige Bailey, Paul Michel, Petko Yotov, Rahma Chaabouni, Ramona 1007 Comanescu, Reena Jana, Rohan Anil, Ross McIlroy, Ruibo Liu, Ryan Mullins, Samuel L Smith, 1008 Sebastian Borgeaud, Sertan Girgin, Sholto Douglas, Shree Pandya, Siamak Shakeri, Soham De, 1009 Ted Klimenko, Tom Hennigan, Vlad Feinberg, Wojciech Stokowiec, Yu hui Chen, Zafarali Ahmed, 1010 Zhitao Gong, Tris Warkentin, Ludovic Peran, Minh Giang, Clément Farabet, Oriol Vinyals, Jeff 1011 Dean, Koray Kavukcuoglu, Demis Hassabis, Zoubin Ghahramani, Douglas Eck, Joelle Barral, 1012 Fernando Pereira, Eli Collins, Armand Joulin, Noah Fiedel, Evan Senter, Alek Andreev, and 1013 Kathleen Kenealy. Gemma: Open models based on gemini research and technology, 2024. URL 1014 https://arxiv.org/abs/2403.08295. 1015
- Adly Templeton, Tom Conerly, Jonathan Marcus, Jack Lindsey, Trenton Bricken, Brian Chen, Adam 1016 Pearce, Craig Citro, Emmanuel Ameisen, Andy Jones, Hoagy Cunningham, Nicholas L Turner, 1017 Callum McDougall, Monte MacDiarmid, C. Daniel Freeman, Theodore R. Sumers, Edward Rees, 1018 Joshua Batson, Adam Jermyn, Shan Carter, Chris Olah, and Tom Henighan. Scaling monoseman-1019 ticity: Extracting interpretable features from claude 3 sonnet. Transformer Circuits Thread, 2024. 1020 URL https://transformer-circuits.pub/2024/scaling-monosemanticity/index.html. 1021
- Hugo Touvron, Louis Martin, Kevin R. Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, Daniel M. Bikel, Lukas 1023 Blecher, Cristian Cantón Ferrer, Moya Chen, Guillem Cucurull, David Esiobu, Jude Fernandes, 1024 Jeremy Fu, Wenyin Fu, Brian Fuller, Cynthia Gao, Vedanuj Goswami, Naman Goyal, Anthony S. 1025 Hartshorn, Saghar Hosseini, Rui Hou, Hakan Inan, Marcin Kardas, Viktor Kerkez, Madian Khabsa,

1049

1077

1026 Isabel M. Kloumann, A. V. Korenev, Punit Singh Koura, Marie-Anne Lachaux, Thibaut Lavril, 1027 Jenya Lee, Diana Liskovich, Yinghai Lu, Yuning Mao, Xavier Martinet, Todor Mihaylov, Pushkar 1028 Mishra, Igor Molybog, Yixin Nie, Andrew Poulton, Jeremy Reizenstein, Rashi Rungta, Kalyan 1029 Saladi, Alan Schelten, Ruan Silva, Eric Michael Smith, R. Subramanian, Xia Tan, Binh Tang, Ross 1030 Taylor, Adina Williams, Jian Xiang Kuan, Puxin Xu, Zhengxu Yan, Iliyan Zarov, Yuchen Zhang, Angela Fan, Melanie Kambadur, Sharan Narang, Aurelien Rodriguez, Robert Stojnic, Sergey 1031 Edunov, and Thomas Scialom. Llama 2: Open foundation and fine-tuned chat models. ArXiv, 1032 abs/2307.09288, 2023. URL https://api.semanticscholar.org/CorpusID:259950998. 1033

- Jesse Vig, Sebastian Gehrmann, Yonatan Belinkov, Sharon Qian, Daniel Nevo, Yaron Singer, and Stuart Shieber. Investigating gender bias in language models using causal mediation analysis. In H. Larochelle, M. Ranzato, R. Hadsell, M.F. Balcan, and H. Lin (eds.), Advances in Neural Information Processing Systems, volume 33, pp. 12388–12401. Curran Associates, Inc., 2020. URL https://proceedings.neurips.cc/paper_files/paper/2020/file/92650b2e92217715fe312e6fa7b90d82-Paper.pdf.
- Kevin Wang, Alexandre Variengien, Arthur Conmy, Buck Shlegeris, and Jacob Steinhardt. Interpretability in the wild: a circuit for indirect object identification in gpt-2 small, 2022. URL https://arxiv.org/abs/2211.00593.
- Jason Wei, Yi Tay, Rishi Bommasani, Colin Raffel, Barret Zoph, Sebastian Borgeaud, Dani Yogatama, Maarten Bosma, Denny Zhou, Donald Metzler, Ed H. Chi, Tatsunori Hashimoto, Oriol Vinyals, Percy Liang, Jeff Dean, and William Fedus. Emergent abilities of large language models. *Transactions on Machine Learning Research*, 2022. ISSN 2835-8856. URL https://openreview.net/forum?id=yzkSU5zdwD. Survey Certification.
- Thomas Wolf, Lysandre Debut, Victor Sanh, Julien Chaumond, Clement Delangue, Anthony Moi, 1050 Pierric Cistac, Tim Rault, Remi Louf, Morgan Funtowicz, Joe Davison, Sam Shleifer, Patrick von 1051 Platen, Clara Ma, Yacine Jernite, Julien Plu, Canwen Xu, Teven Le Scao, Sylvain Gugger, Mariama 1052 Drame, Quentin Lhoest, and Alexander Rush. Transformers: State-of-the-art natural language 1053 processing. In Qun Liu and David Schlangen (eds.), Proceedings of the 2020 Conference on 1054 Empirical Methods in Natural Language Processing: System Demonstrations, pp. 38–45, Online, 1055 October 2020. Association for Computational Linguistics. doi: 10.18653/v1/2020.emnlp-demos.6. 1056 URL https://aclanthology.org/2020.emnlp-demos.6. 1057
- Bingyang Wu, Ruidong Zhu, Zili Zhang, Peng Sun, Xuanzhe Liu, and Xin Jin. {dLoRA}: Dynamically orchestrating requests and adapters for {LoRA}{LLM} serving. In *18th USENIX Symposium on Operating Systems Design and Implementation (OSDI 24)*, pp. 911–927, 2024a.
- Zhengxuan Wu, Atticus Geiger, Thomas Icard, Christopher Potts, and Noah Goodman. Interpretability at scale: Identifying causal mechanisms in alpaca. In A. Oh, T. Naumann, A. Globerson, K. Saenko, M. Hardt, and S. Levine (eds.), *Advances in Neural Information Processing Systems*, volume 36, pp. 78205–78226. Curran Associates, Inc., 2023. URL https://proceedings.neurips.cc/paper_files/paper/2023/file/f6a8b109d4d4fd64c75e94aaf85d9697-Paper-Conference.pdf.
- Zhengxuan Wu, Atticus Geiger, Aryaman Arora, Jing Huang, Zheng Wang, Noah Goodman, Christopher Manning, and Christopher Potts. pyvene: A library for understanding and improving PyTorch models via interventions. In Kai-Wei Chang, Annie Lee, and Nazneen Rajani (eds.), *Proceedings of the 2024 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 3: System Demonstrations)*, pp. 158–165, Mexico City, Mexico, June 2024b. Association for Computational Linguistics. URL https://aclanthology.org/2024.naacl-demo.16.
- An Yang, Baosong Yang, Binyuan Hui, Bo Zheng, Bowen Yu, Chang Zhou, Chengpeng Li, Chengyuan Li, Dayiheng Liu, Fei Huang, et al. Qwen2 technical report. arXiv preprint arXiv:2407.10671, 2024.
- Alex Young, Bei Chen, Chao Li, Chengen Huang, Ge Zhang, Guanwei Zhang, Heng Li, Jiangcheng Zhu, Jianqun Chen, Jing Chang, et al. Yi: Open foundation models by 01. ai. arXiv preprint arXiv:2403.04652, 2024.

 Fred Zhang and Neel Nanda. Towards best practices of activation patching in language models: Metrics and methods. In *The Twelfth International Conference on Learning Representations*, 2024. URL https://openreview.net/forum?id=Hf17y6u9BC.

1084	Susan Zhang, Stephen Roller, Naman Goyal, Mikel Artetxe, Moya Chen, Shuohui Chen, Christopher
1085	Dewan, Mona Diab, Xian Li, Xi Victoria Lin, Todor Mihaylov, Myle Ott, Sam Shleifer, Kurt
1086	Shuster, Daniel Simig, Punit Singh Koura, Anjali Sridhar, Tianlu Wang, and Luke Zettlemoyer.
1087	Opt: Open pre-trained transformer language models, 2022. URL https://arxiv.org/abs/2205.
1088	01068.
1089	
1090	
1091	
1092	
1093	
1094	
1095	
1096	
1097	
1098	
1099	
1100	
1101	
1102	
1103	
1104	
1105	
1106	
1107	
1108	
1109	
1110	
1111	
1112	
1113	
1114	
1115	
1116	
1117	
1118	
1119	
1120	
1121	
1122	
1123	
1124	
1125	
1120	
1127	
1120	
1129	
1121	
1132	
1122	
1100	

1134 A RESEARCH SURVEY DETAILS

In this section, we provide details on the data collection process for the results reported in Section 2, and report additional findings related to model parameter size.

Data Curation. We curated a list of 184 total papers that study the internals of open-weight transformer models, taken from the citations of a recent survey of interpretability research (Ferrando et al., 2024). Starting from the initial list of 411 citations, we excluded any paper that met one or more of the following criteria:

- Paper is a survey paper, tutorial, library report, model card, or performs no model experiments.
- Paper only performs experiments on closed-weight models, custom models, or non-transformer architectures.
- Paper was published before 2019.

We identified the largest model (by parameter count) studied in each paper and gathered data related to its MMLU performance (Hendrycks et al., 2021) from the HuggingFace OpenLLM Leaderboard Archive (Huggingface, 2024). We used the Papers with Code MMLU Leaderboard (PaperswithCode, 2024) to supplement MMLU performance data when models did not have a score on the Huggingface leaderboard. When a model had a score from both sources, we averaged its benchmark results together. A few models did not have reported MMLU results from either source (primarily BERT-style models). For these models, we interpolated MMLU results using nearest neighbor approximation based on parameter size. Papers with Code provides an MMLU performance reference for closed-weight models.



Figure 7: Evolution of the size gap between open-weight models used in research and publicly released models from 2019 to 2024. Blue boxes represent the distribution of model sizes used in research papers, while pink boxes show publicly released models. The dashed gold line connects median model sizes for each group, with the ratio between medians displayed. Outliers are shown as individual points. The growing disparity (from 2.7x in 2019-2020 to 10.3x in 2024) suggests that interpretability research is increasingly lagging behind the capabilities of state-of-the-art models available to the public.

Comparing Model Parameter Size. In Figure 7, we show there is a growing disparity between the size of models that interpretability researchers are studying and the size of models that are publicly released each year, a similar trend to the one discussed in §2 regarding MMLU. Reference data on parameter sizes of publicly available models is taken from Epoch AI (2024), restricting our analysis to language models released after 2019 that have openly downloadable parameters. To ensure reproducibility and transparency, all source code and data used for the analyses described in this section are available in the supplementary materials accompanying this paper.

¹¹⁸⁸ B IMPLEMENTATION DETAILS

1189 1190

NNsight and Envoys To support the tracing context, the NNsight library wraps PyTorch modules 1191 in an NNsight object instance. The NNsight wrapper is responsible for providing access to its sub-1192 modules' intermediate deferred inputs and outputs. Upon initialization, the NNsight object creates 1193 an Envoy object for each sub-module, forming a tree-like structure mirroring that of the original 1194 PyTorch module tree. Each Envoy is responsible for managing and recording operations on future 1195 inputs and outputs for its underlying module. It achieves this by exposing attributes like .input and .output which, when accessed, trace a deferred operation to intervene at the corresponding 1196 module's input or output during model execution. This attribute access returns a Proxy object of 1197 that deferred operation, representing a future intermediate value. Any operation performed on the 1198 resulting Proxy object creates a new deferred operation, and therefore a new Proxy. As all future 1199 Proxys originate from operations performed on these root . input and .output attributes, it makes 1200 them the entry-point into the tracing context functionality. 1201

```
1202 1 from nnsight import NNsight
```

```
1203 2 model = NNsight(model)
```

```
1204 3 input = tensor([.46, 1.56, -0.3, .98, -0.5])
1205 4 with model.trace(input):
1206 5 model.submodule.output[:2] = 0
1206 6 hidden_state = model.submodule.output.save()
```

1207 Code Example 3: Basic usage of the NNsight tracing context and the NNsight data type. model 1208 in this example has a single submodule named "submodule" and accepts a Tensor of length five. On 1209 line two, the PyTorch model is wrapped with NNsight. On line four, the tracing context is created 1210 and entered by calling with model.trace(...) given some input. On line five, the Envoy for 1211 "submodule" is used to access a Proxy for its deferred intermediate output via the .output attribute. 1212 Also on line five, the returned Proxy is intervened on, creating two new Proxys and Nodes through 1213 the overloaded slicing and setting methods. Finally on line six, the same intermediate value is placed 1214 in the "hidden_state" variable, and saved for use after model execution.

1215

When exiting the tracing context, the NNsight object interleaves the traced operations to the model's computation graph and executes the model. This process is performed by adding PyTorch hooks to all modules whose Envoy's .input or .output attributes were accessed within the tracing context. These hooks are the interface between the intervention graph, formed from the deferred operations gathered while inside of the tracing context, and the model's original computation graph. After interleaving, NNsight removes these hooks.

1222

1223 **The Intervention Graph** An intervention graph is composed of Node objects. Each Node represents 1224 a single operation, or target, to be executed during evaluation of the graph. Nodes have a unique 1225 identifier, or name, and arbitrary positional and keyword arguments. These arguments may contain 1226 other Nodes which are considered *dependencies* of the source Node. Conversely, the source Node 1227 is a *listener* of its *dependency*. When a Node gets executed during interleaving, it decrements the 1228 amount of remaining *dependencies* of its *listeners* and decrements the amount of remaining *listeners* of its *dependencies*. If the amount of remaining *dependencies* of a Node becomes zero, then that 1229 Node is immediately executed. This means that the root intervention Nodes created with Envoys are 1230 the only Nodes executed outside of this dependency chain, and setting them through interleaving 1231 "kicks off" the rest of the Intervention Graph. When a Node's remaining listeners reaches zero, it 1232 is considered redundant and has its value is cleared and its memory freed immediately. 1233

For every Node, there is a corresponding Proxy that manages tracing operations on the underlying Node. Proxys achieve this by overriding builtin Python and PyTorch magic methods. These overrides intercept calls to methods and functions that now create and add a new Node to the Intervention Graph, returning the Node's corresponding Proxy. Chaining this process results in a completed Intervention Graph.

1239

Protocols When composed into an Intervention Graph, Nodes are quite powerful on their own.
 However, a full-featured framework requires operations that can exert more control on the intervention process than the scope an individual Node can provide.

1242 PyTorch **Research Code Intervention Graph** 1243 LLama from nnsight import LanguageModel InterventionProtocolNode(Im.output, FakeTensor(data)) 1244 Model key = 'meta-llama/Meta-Llama-3-8B' LockProtocolNode() 1245 lm = LanguageModel(key) Envoy 1246 with lm.trace('Hello World',scan=True): If scan=True, populate Nodes with FakeTensor(...) value = lm.output 1247 value.save() by using module shape values in Pytorch LLama model Proxy

1248

Figure 8: The components involved in creating nodes within the Intervention Graph. Module calls like lm.output are processed through Envoy, resulting in the creation of relevant ProtocolNodes like InterventionProtocolNode(). Everything else goes through Proxy layers to create other relevant nodes in the graph. All nodes are encapsulated within InterventionProxy(), which users see upon printing nodes. NNsight can scan the modules in the PyTorch model instance to gather related data shapes with a scan flag, creating FakeTensors that are included in the respective intervention graph nodes and returned to the user.

1256

1257

A Protocol is a special NNsight operation that enables these more powerful controls. A Protocol 1259 must define two methods: One that adds a Node of its type to the Intervention Graph, and another 1260 that executes that Node during interleaving. Rather than NNsight executing a Node's target with its 1261 arguments, Protocols allow the Node to be passed the Protocol for execution. With full access to the Node, the Protocol can processes its arguments, implement complex logic, and access the 1262 Intervention Graph itself. This access provides the Protocol a global view of the intervention 1263 process, which it can leverage by interacting with the attachments attribute on the Intervention 1264 Graph. Attachments are a dictionary on the Intervention Graph where Protocols can manage a 1265 global state throughout interleaving. 1266

1267 Two types of Protocols are highlighted below:

LockProtocol: When a Node has no more listeners, its value is cleared to free memory. Therefore this NO-OP operation never updates the remaining listeners of its dependencies, preventing a Node's value from being deleted. Proxys implement a .save() method which creates a LockProtocol with the source Node as its only dependency, resulting in the Node's value being available outside the tracing context.

GradProtocol: Many use cases require intervening on the backwards pass of a model rather than exclusively the forward pass. Accessing a Proxy's.grad attribute creates a GradProtocol Node.
When executed, GradProtocol adds a backwards hook to the underlying Tensor value. This hook injects the gradients into the GradProtocol Node during the backwards pass, allowing it to behave like any other Node.

1278

Model Abstractions The NNsight class defines abstract convenience methods that subclasses can implement to extend common neural network architectures to the NNsight API. These methods include logic to load the model, which is called twice when initializing and executing an NNsight-based model from scratch. The first call happens upon initialization to create and populate the Envoy tree with a 'meta' version of the model. In a 'meta' model, the parameters of the network are not yet loaded: they exist as placeholders containing shape and datatype information. The second call occurs when loading and dispatching the true model parameters.

1286 NNsight comes with a few of these subclasses pre-implemented. The most common subclass 1287 is the LanguageModel, for text-to-text sequence generation architectures. The LanguageModel 1288 integrates with Huggingface to load model configurations, tokenizers, and parameters from their 1289 online repository when provided with a repo id. Implementing the input processing methods uses the model's tokenizer to accept common forms of inputs and batch them together. When loading the 1290 'meta' version of the model, the LanguageModel uses AutoModelForCausalLM.from_config from 1291 Huggingface transformers to download only the models configuration. Only upon dispatching 1292 does LanguageModel fetch the model parameters if they aren't already cached. 1293

- 1294
- **Scanning and Validation** NNsight offers two key debugging features integrated with the Tracer context.

The first is Scanning. When toggled, scanning performs an initial execution pass on the model in
 FakeTensorMode. During this pass, Hooks registered on each module update the input and output of
 their corresponding Envoys on the NNsight model with generated FakeTensors. This process makes
 valuable information available during tracing such as output shapes, dtypes, and device allocations.

When Scanning is enabled, Validation can be used to test interventions on FakeTensors as they are being added on the intervention graph. This acts as a sanity check to identify errors related to defined interventions before engaging in the model's forward pass run.

Remote Execution and Session It is common for users to run multiple forward passes on a model as part of a larger experiment, such as in a LORA training loop. NNsight offers an API that allows the creation of a Session context, where inner defined Tracer contexts are executed sequentially only after the Session context is exited. By encapsulating multiple Tracers in a Session, values obtained in earlier passes can be seamlessly referenced by later stages in the experiment without the need for manual saving.

When a Session is executed remotely, all forward passes defined within are sent as part of a single request to be executed. This can reduce wait times between Tracer calls by 1) minimizing the number of server requests and 2) eliminating the need to download and re-upload values referenced across Tracer contexts for subsequent calls that require them.

¹³⁵⁰ C CODE COMPARISON

1351

In Section 3.2, we discussed how NNsight can express any experiment that can be expressed with PyTorch, but that is also simplifies the design of custom experiments on large model internals compared to using standard PyTorch hooks. Here we provide an additional pair of code examples that implement a common interpretability technique called "activation patching" (Zhang & Nanda, 2024). Code Example 4 implements this experiment using PyTorch hooks, while Code Example 5 uses the NNsight API. Both code snippets implement the same intervention, causing the model to change its prediction for the base prompt from "Paris" to "Rome".

```
from transformers import AutoTokenizer, AutoModelForCausalLM
model_id = "meta-llama/Meta-Llama-3.1-8B"
1359
         2
1360
         3 tokenizer = AutoTokenizer.from_pretrained(model_id, padding_side='left')
         4 tokenizer.pad_token = tokenizer.eos_token
         5 lm = AutoModelForCausalLM.from_pretrained(model_id)
1362
1363
         7 edit_prompt = "The Colosseum is located in the city of"
         8 base_prompt = "The Louvre is located in the city of'
1364
         9 batch = [edit_prompt,base_prompt]
1365
        11 edit tok = 5
1366
        12 base_tok = 6
1367
        14 def hook_fn(module, input, output):
1368
               output[0][1,base_tok,:] = output[0][0,edit_tok,:]
        15
1369
        17 layer = lm.model.layers[5]
1370
        18 hook = layer.register_forward_hook(hook_fn)
1371
        19 inputs = tokenizer([edit_prompt,base_prompt], return_tensors="pt", padding=True)
        20 output = lm(**inputs)
1372
        21 hook.remove()
1373
        23 last = output["logits"][-1,-1].argmax()
1374
        24 prediction = tokenizer.decode(last)
1375
        25 print(prediction)
```

Code Example 4: An intervention implemented with standard PyTorch hooks. Lines 1-4 load a language model from huggingface, and lines 7-9 define two prompts. In lines 11-12, we note the token indices for the last token of the subjects in each prompt. Lines 14-15 define a PyTorch hook that replaces the hidden state in the base prompt with the hidden state from the edit prompt using the respective last-subject-token indices. The experiment is carried out in lines 17-21, and we collect and print the output of the model in lines 23-25.

```
1 from nnsight import LanguageModel
         2 model_id = "meta-llama/Meta-Llama-3.1-8B"
1384
        3 lm = LanguageModel(model_id)
1385
         5 edit_prompt = "The Colosseum is located in the city of"
1386
         6 base_prompt = "The Louvre is located in the city of
        7 batch = [edit_prompt,base_prompt]
1387
1388
        9 edit_tok = 5
1389
        10 base_tok = 6
1390
        12 layer = lm.model.layers[5]
1391
        13 with lm.trace(batch) as tracer:
        14
              layer.output[0][1,base_tok,:] = layer.output[0][0,edit_tok,:]
1392
        15
               output = lm.output.save()
1393
        17 last = output["logits"][-1,-1].argmax()
1394
        18 prediction = lm.tokenizer.decode(last)
1395
        19 print(prediction)
```

Code Example 5: An intervention implemented with the NNsight API. This code defines the same intervention as in Code Example 4. Lines 1-3 load a language model from huggingface, and lines 5-7 define two prompts. In lines 9-10, we note the token indices for the last token of the subjects in each prompt. The experiment is carried out in lines 12-15, where we replace the hidden state in the base prompt with the hidden state from the edit prompt, using the respective last-subject-token indices. We collect and print the output of the model in lines 17-19.

1402

1382

¹⁴⁰³ In Code Example 6, we show an example of how to implement and train a LoRA (Hu et al., 2021) adapter using a session context with remote execution via NDIF. Code Examples 7 and 8 demonstrate

interventions on the gradients of model computations. And in Code Example 9, we show an what
 training a linear probe remotely via NDIF might look like to predict output of layer 1 using the output of layer 0.

```
1 import torch
1408
1409
         3 from nnsight.envoy import Envoy
         4 from torch.utils.data import DataLoader
1410
1411
         6 # We will define a LORA class.
         7 \# The LORA class call method operations are simply traced like you would normally do in a .
1412
                trace
1413
         8 class LORA(nn.Module):
               def __init__(self, module: Envoy, dim: int, r: int) -> None:
    """Init.
1414
         9
        10
1415
1416
        12
                   Args:
        13
                        module (Envoy): Which model Module we are adding the LORA to.
1417
        14
                       dim (int): Dimension of the layer we are adding to (This could potentially be
1418
                auto populated if the user scanned first so we know the shape)
        15
                        r (int): Inner dimension of the LORA
1419
        16
                   super(LORA, self).__init__()
1420
        18
                   self.r = r
1421
        19
                   self.module = module
                   self.WA = torch.nn.Parameter(torch.randn(dim, self.r), requires_grad=True).save()
1422
        20
                   self.WB = torch.nn.Parameter(torch.zeros(self.r, dim), requires_grad=True).save()
        21
1423
        23
               # The Call method defines how to actually apply the LORA.
1424
               def __call__(self, alpha: float = 1.0):
    """Call.
        24
1425
        25
1426
        27
                   Args:
1427
                       alpha (float, optional): How much to apply the LORA. Can be altered after
        28
                training for inference. Defaults to 1.0.
1428
        29
1429
                   # We apply WA to the first positional arg (the hidden states)
1430
        31
                   A_x = torch.matmul(self.module.input[0][0], self.WA)
1431
                   BA_x = torch.matmul(A_x, self.WB)
        33
1432
        35
                   # LORA is additive
1433
        36
                   h = BA_x + self.module.output
1434
        38
                   # Replace the output with our new one * alpha
1435
        39
                   # Could also have been self.module.output[:] = h * alpha, for in-place
        40
                   self.module.output = h * alpha
1436
1437
        42
               def parameters(self):
                   # Some way to get all the parameters.
1438
        43
                   return [self.WA, self.WB]
        44
1439
        46 # get the model
1440
        47 llama = LanguageModel("meta-llama/Meta-Llama-3.1-70B")
1441
        49 # We need the token id of the correct answer.
50 answer = " Paris"
1442
1443
        51 answer_token = llama.tokenizer.encode(answer)[1]
        52 # Inner LORA dimension
1444
        53 \text{ lora}_\text{dim} = 4
1445
        55 # Module to train LORA on
1446
        56 module = llama.model.layers[-1].mlp
1447
1448
1449
        _{60} # The LORA object itself isn't transmitted to the server. Only the forward / call method.
        61 # The parameters are created remotely and never sent only retrieved
1450
        62 with llama.session(remote=True) as session:
1451
               # Create dataset of 100 pairs of a blank prompt and the " Paris " id
1452
        64
               dataset = [["_", answer_token]] * 100
        65
1453
        67
               # Create a dataloader from it.
1454
        68
               dataloader = DataLoader(dataset, batch_size=10)
1455
        70
               # Create our LORA on the last mlp
1456
               lora = LORA(module, dim, lora_dim)
1457
        73
               # Create an optimizer. Use the parameters from LORA
```

```
1458
         74
                optimizer = torch.optim.AdamW(lora.parameters(), lr=3)
1459
         76
                # Iterate over dataloader using .iter.
1460
                with session.iter(dataloader, return_context=True) as (batch, iterator):
         77
1461
         79
                     prompt = batch[0]
1462
         80
                     correct_token = batch[1]
1463
                     # Run .trace with prompt
         82
1464
                    with llama.trace(prompt) as tracer:
         83
1465
                         # Apply LORA to intervention graph just by calling it with .trace
1466
         85
         86
                         lora()
1467
         88
                         # Get logits
1468
                        logits = llama.lm head.output
         89
1469
                         # Do cross entropy on last predicted token and correct_token
         91
1470
                         loss = torch.nn.functional.cross_entropy(logits[:, -1], batch[1])
         92
1471
                         # Call backward
         93
        94
                         loss.backward()
1472
1473
                    # Call methods on optimizer. Graphs that arent from .trace (so in this case session
        96
                 and iterator both have their own graph) are executed sequentially.
1474
        97
                   # The Graph of Iterator here will be:
1475
                    # 1.) Index batch at 0 for prompt
         98
                    # 2.) Index batch at 1 for correct_token
         99
1476
        100
                    # 3.) Execute the .trace using the prompt
1477
        101
                    # 4.) Call .step() on optimizer
        102
                    optimizer.step()
1478
        103
                     # 5.) Call .zero_grad() in optimizer
1479
        104
                     optimizer.zero_grad()
        105
                     # 6.) Print out the lora WA weights to show they are indeed changing
1480
                    iterator.log(lora.WA)
        106
1481
         Code Example 6: Using NNsight to train a LORA with remote execution, the Session context, and
1482
         iterative interventions. First, the LORA class is defined that takes an Envoy as input. Then, the
1483
         session context is used for remote training with an optimizer.
1484
1485
         1 import nnsight
1486
         2 from nnsight import NNsight
1487
         3 from collections import OrderedDict
          4 import torch
1488
1489
         6 input size = 5
          7 hidden_dims = 10
1490
         8 output_size = 2
1491
         10 net = torch.nn.Sequential(
1492
              OrderedDict(
        11
1493
         12
                    Ε
                         ("layer1", torch.nn.Linear(input_size, hidden_dims)),
("layer2", torch.nn.Linear(hidden_dims, output_size)),
1494
         14
1495
                    ٦
         15
                )
         16
1496
        17 ).requires_grad_(False)
18 tiny_model = NNsight(net)
1497
        19 input = torch.rand((1, input_size))
1498
1499
        21 with tiny_model.trace(input):
1500
                # We need to explicitly have the tensor require grad
         23
                # as the model we defined earlier turned off requiring grad.
1501
         24
         25
                tiny_model.layer1.output.requires_grad = True
1502
1503
         27
                \ensuremath{\texttt{\#}} We call .grad on a tensor Proxy to communicate we want to store its gradient.
                # We need to call .save() since .grad is its own Proxy.
layer1_output_grad = tiny_model.layer1.output.grad.save()
         28
1504
         29
1505
                layer2_output_grad = tiny_model.layer2.output.grad.save()
         30
1506
         32
                # Need a loss to propagate through the later modules in order to have a grad.
1507
         33
                loss = tiny_model.output.sum()
        34
                loss.backward()
1508
        36 print("Layer 1 output gradient:", layer1_output_grad)
37 print("Layer 2 output gradient:", layer2_output_grad)
1509
1510
```

1511 Code Example 7: Applying backpropagation and accessing gradients with respect to a loss.

```
1512
         1 import nnsight
1513
         2 from nnsight import NNsight
         3 from collections import OrderedDict
1514
         4 import torch
1515
         6 input_size = 5
1516
         7 hidden_dims = 10
1517
         8 output_size = 2
1518
        10 net = torch.nn.Sequential(
1519
               OrderedDict(
        11
                  Ε
1520
                         ("layer1", torch.nn.Linear(input_size, hidden_dims)),
("layer2", torch.nn.Linear(hidden_dims, output_size)),
         13
1521
         14
         15
                    J
1522
               )
         16
1523
         17 ).requires_grad_(False)
         18 tiny_model = NNsight(net)
1524
        19 input = torch.rand((1, input_size))
1525
        21 with tiny_model.trace(input):
1526
1527
        23
                # We need to explicitly have the tensor require grad
                # as the model we defined earlier turned off requiring grad.
         24
1528
        25
                tiny_model.layer1.output.requires_grad = True
1529
         27
                tiny_model.layer1.output.grad[:] = 0
1530
        28
                tiny_model.layer2.output.grad = tiny_model.layer2.output.grad * 2
1531
        30
                layer1_output_grad = tiny_model.layer1.output.grad.save()
1532
        31
                layer2_output_grad = tiny_model.layer2.output.grad.save()
1533
        33
                \ensuremath{\texttt{\#}} Need a loss to propagate through the later modules in order to have a grad.
1534
         34
                loss = tiny_model.output.sum()
1535
        35
                loss.backward()
1536
        37 print("Layer 1 output gradient:", layer1_output_grad)
38 print("Layer 2 output gradient:", layer2_output_grad)
1537
1538
            Code Example 8: Zeroing out the grad of one layer, and doubling the grad of a second layer.
1539
1540
         1 import torch
1541
         2 from torch.utils.data import DataLoader
         3 from nnsight import LanguageModel, list, log
1542
1543
         5 model = LanguageModel("meta-llama/Meta-llama-3.1-405B")
1544
         7 with model.session(remote=True) as session:
1545
                features = 4096
         8
         9
                probe = torch.nn.Linear(features, features).save()
1546
         10
                dataset = DataLoader(
                                       "to", "train", "on"], batch_size=2, shuffle=True
1547
        11
                    ["some", "text",
         12
                )
1548
         13
                optimizer = torch.optim.Adam(probe.parameters(), lr=0.003)
1549
         15
                for epoch in list(range(50)):
1550
         16
                     for batch in dataset:
1551
                         with model.trace(batch) as tracer:
         18
                             input = model.model.layers[0].output[0]
1552
         19
                             probe = probe.to(input.device)
1553
        20
                             prediction = probe(input)
                             output = model.model.layers[1].output[0]
        21
1554
                             loss = torch.nn.functional.mse_loss(prediction, output)
1555
                             loss.backward()
1556
         25
                         optimizer.step()
1557
         26
                         optimizer.zero_grad()
1558
        28
                    log(loss)
1559
        29
                    log(probe.weight)
1560
              Code Example 9: Sample code showing the remote training of a linear probe using NDIF.
1561
1562
1563
1564
1565
```

1566 D PERFORMANCE

In this section, we report additional results evaluating the performance of NNsight and NDIF in various settings. These include setup and runtime across different model sizes of OPT (Zhang et al., 2022), comparing the setup and runtime of HPC and NDIF for Llama-3.1 models (Dubey et al., 2024), and the response time for NDIF requests as a function of the number of users.

1573 D.1 HPC vs. NDIF

Table 2: Setup-time and runtime comparison of HPC and NDIF across OPT parameter sizes

OPT	Parameters	Setup Time	Runtime
	125m	0.635 ± 0.035	0.021 ± 0.001
	350m	0.847 ± 0.058	0.041 ± 0.001
	1.3b	1.242 ± 0.016	0.072 ± 0.004
HPC	2.7b	1.866 ± 0.003	0.123 ± 0.011
	6.7b	3.474 ± 0.023	0.257 ± 0.008
	13b	5.971 ± 0.021	0.429 ± 0.018
	30b	11.338 ± 0.073	0.913 ± 0.016
	66b	23.697 ± 0.079	1.889 ± 0.040
	125m	0.149 ± 0.008	0.541 ± 0.123
	350m	0.304 ± 0.015	0.588 ± 0.178
	1.3b	0.409 ± 0.213	0.627 ± 0.099
NDIF	2.7b	0.358 ± 0.054	0.494 ± 0.085
	6.7b	0.317 ± 0.011	0.550 ± 0.080
	13b	0.370 ± 0.026	0.773 ± 0.136
	30b	0.343 ± 0.013	1.295 ± 0.104
	66b	0.423 ± 0.095	2.739 ± 0.154

Table 3: Runtime comparison of NNsight using execution on an HPC and remote execution via NDIF.

Framework	Activation Patching		
	Llama-3.1-8B	Llama-3.1-70B	
NNsight (HPC) NNsight (NDIF)	$\begin{array}{c} 0.333 \pm 0.013 \\ 0.998 \pm 0.134 \end{array}$	$\begin{array}{c} 1.992 \pm 0.056 \\ 2.369 \pm 0.117 \end{array}$	

Table 4: Time in seconds to load the model into memory.

Framework	Loading Weights		
	Llama-3.1-8B	Llama-3.1-70B	
NNsight (HPC) NNsight (NDIF)	$\begin{array}{c} 5.991 \pm 0.192 \\ 0.518 \pm 0.005 \end{array}$	$\begin{array}{c} 43.617 \pm 1.118 \\ 0.695 \pm 0.056 \end{array}$	

1614 D.2 LOAD-TESTING NDIF FOR CONCURRENT USERS

In this section, we assess the performance and scalability of NDIF with an increasing number of concurrent users (and thereby requests). We use Llama-3.1-8B (Dubey et al., 2024) as the model being served on NDIF for this evaluation, which was hosted on a 48 GB RTX 6000 Ada GPU.

1619 To simulate a single user, we construct a NNsight request with a prompt containing up to 24 tokens that accesses and saves the output of a layer of Llama-3.1-8B selected uniformly at random. Code

Example 10 shows how we generate sample prompts and requests for a single user, saving the response time of the request sent to and returned by NDIF.

```
def layer_selector(self) -> None:
1623
                """Simulate selecting activations of an intermediate layer.""" start_time = time.time()
1624
         4
1625
                    layer : int = random.randint(0, self.n_layers - 1)
                    query : str = 'hello ' * random.randint(1, 24)
1626
         6
                    with self.model.trace(query, remote=True):
1627
                        output = self.model.model.layers[layer].output.save()
         9
                finally:
1628
        10
                    elapsed_time = time.time() - start_time
1629
                    # Append to a list visible locally to this process
                    self.user.environment.request_times.append(elapsed_time)
1630
```

Code Example 10: Code to generate NDIF requests and measure the response time for many users.

Then, for $N \in \{1, ..., 100\}$, we simulate N concurrent users submitting requests to the NDIF server, all running as separate processes, and record the response time for each user's request.

The results can be seen in Figure 9, where we plot the response time of the NDIF system as a function of the number of concurrent users. We find that as the number of users increases, the increase in the median response time (shown in black) is approximately linear. Furthermore, as more concurrent users access the system, variance in response time increases.

This experiment was done using an implementation of the system that creates a queue for each subsequent user, and runs multiple forward passes (one for each user). There are future plans to allow true parallel execution of user requests in the same batch.



Figure 9: Speed and robustness of NDIF response time to changes in concurrent user request count. Tests simulated N = 1 to 100 concurrent users on Llama-3.1-8B requesting outputs from uniformly random intermediate layers over a 30-second window. The gray area shows 25% and 75% quantiles. Points show the minimum and maximum times for each user count. Line plot shows median response time.

1665 1666 1667

1668

1661

1662

1663

1664

1631 1632

- 1669
- 1670
- 1671
- 1672
- 1673

1675 E EQUIVALENCE OF THEANO'S MANY-TO-MANY AND NNSIGHT'S MANY-TO-ONE APPLY NODES

We demonstrate that Al-Rfou et al. (2016)'s many-to-many computation graph structure—where an apply node can output multiple variable nodes—is computationally equivalent to a graph in which each apply node outputs only a single variable node. This equivalence holds by considering the many-to-many structure as a special case of a many-to-one structure, where the single output is a tuple (or concatenation) of variable nodes.

Let $a^{(in)}$ denote an apply node in the many-to-many setting, and let v_1, \ldots, v_m represent the m variable node outputs of $a^{(in)}$. Suppose these variable nodes serve as inputs to n downstream apply nodes $a_1^{(out)}, \ldots, a_n^{(out)}$ in the computation graph.

We now construct a new graph that adheres to the many-to-one structure, where each apply node has exactly one output. Let $b^{(in)}$, v, and $b_1^{(out)}$,..., $b_n^{(out)}$ denote corresponding nodes in the transformed graph, defined as follows:

- 1. The apply node $b^{(in)}$ outputs a single variable node v, which is defined as the concatenation (or tuple) of the original outputs: $v = (v_1, \ldots, v_m)$.
- 2. Each downstream apply node $b_i^{(\text{out})}$ takes as input v and any other variables that were inputs to $a_i^{(\text{out})}$ in the original graph. The node $b_i^{(\text{out})}$ performs the same computation as $a_i^{(\text{out})}$, but accesses the necessary components of v corresponding to its original inputs. Specifically, if $a_i^{(\text{out})}$ originally took inputs v_{i_1}, \ldots, v_{i_k} , then $a_i^{(\text{out})}$ extracts v_{i_1}, \ldots, v_{i_k} from v and processes them accordingly.

This process generalizes inductively to any depth of the computation graph. In the event that one of the output variable nodes v_i needs to be used elsewhere in the computation graph (e.g., as input to a node that does not directly accept v), the transformation can be handled by introducing an additional apply node, $b^{(\text{extract})}$. This node takes v as input and outputs v_i , ensuring that the structure of the original computation graph is maintained.