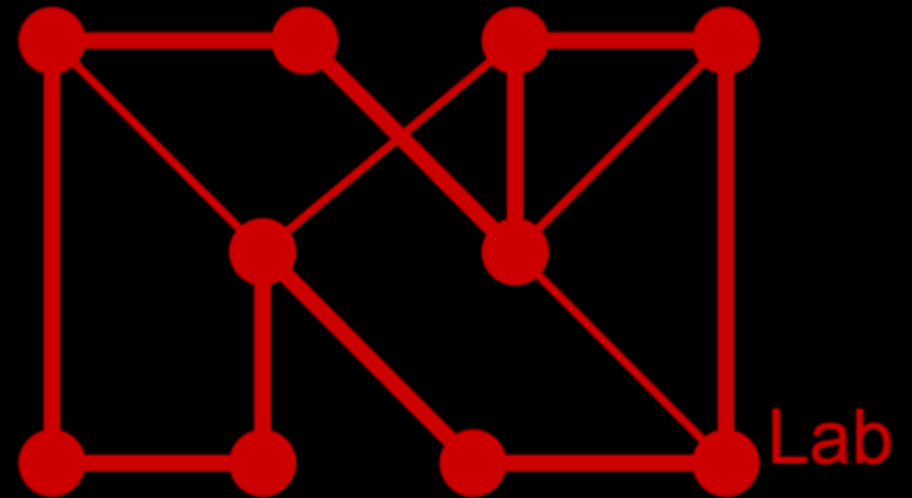


NNSIGHT AND NDIF: DEMOCRATIZING ACCESS TO OPEN-WEIGHT FOUNDATION MODEL INTERNALS

NEURAI Lab, Silicon Valley

Jose L Sampedro Mazon
sampedromazon.j@northeastern.edu



Dario Amodei The Urgency of Interpretability

Dario Amodei — The Urgency of Interpretability



ANTHROPIC

“AI systems are grown more than they are built—their internal mechanisms are “emergent” rather than directly designed (...) we have no idea, at a specific or precise level, why it makes the choices it does, or why it occasionally makes a mistake despite usually being accurate.

(...)

I am very concerned about deploying such systems without a better handle on interpretability. These systems will be absolutely central to the economy, technology, and national security, and will be capable of so much autonomy that I consider it basically unacceptable for humanity to be totally ignorant of how they work.”

AI models offer different levels of access

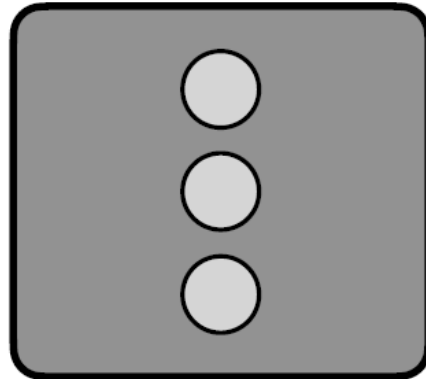
Black-box access allows users to design inputs for a system, query it and analyze the *outputs* (e.g. *OpenAI ChatGPT / completions API*).

Black-Box



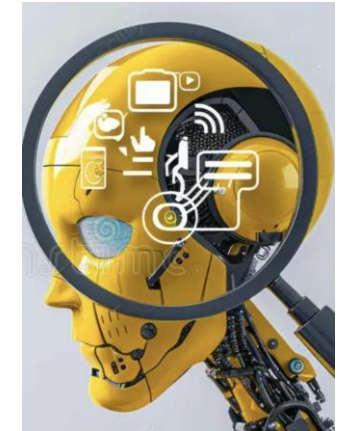
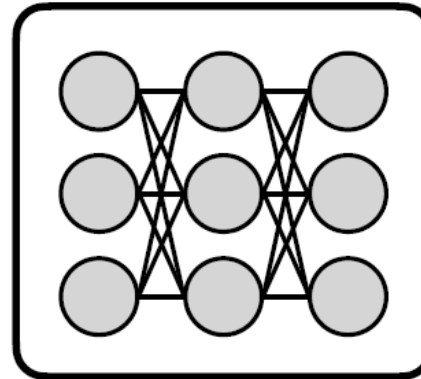
Grey-box models allow access to intentionally exposed model internals like inner neuron activations, input embeddings, probabilities, etc. (e.g. Cohere)

Grey-Box

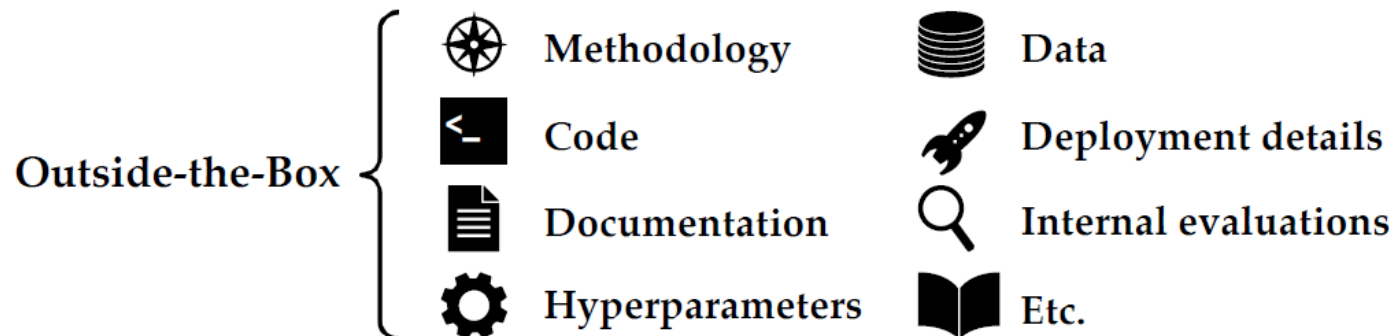


'De facto' white-box access allows running custom tests programmatically via API but system's parameters cannot be copied (e.g. PyTorch Hub models)

White-Box



White-box access allows users full access to weights, activations, gradients, and enable model fine-tuning (e.g. *Llama 3, Falcon*)



Outside-the-box access can include methodological details, source code, documentation, hyper-parameters, training data, deployment details, and findings from internal evaluations (e.g. *Hugging Face Model Cards*)



Level of access dictates depth of interpretability studies

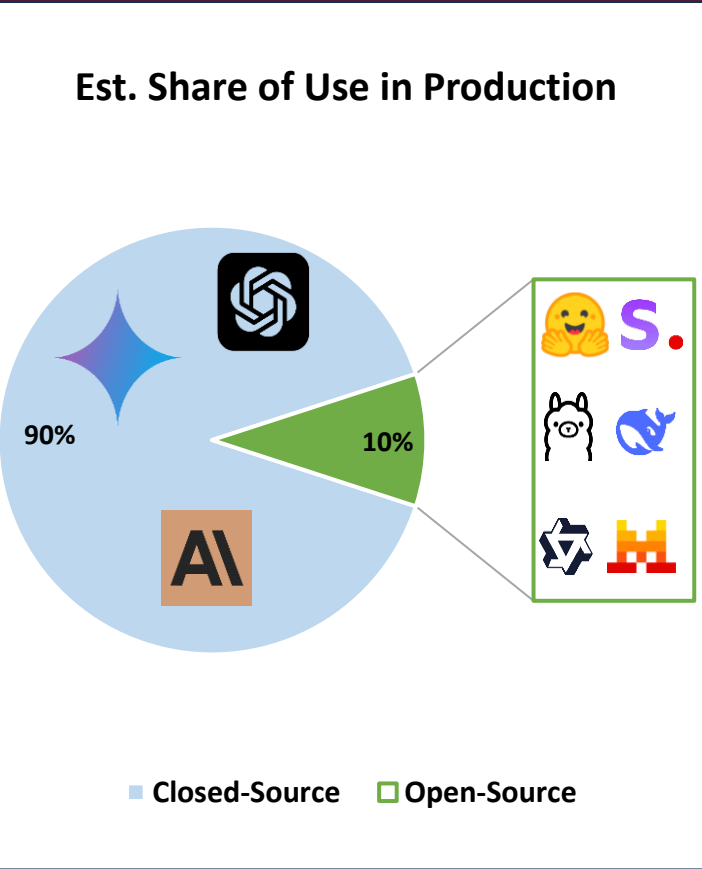
Access Type	Interpretability Potential	Key Limitations
Black-Box <i>(Prompt system and observe output only)</i>	<ul style="list-style-type: none">✓ Test inputs & outputs✓ Assess basic performance✓ Identify obvious biases	<ul style="list-style-type: none">✗ Hard detect rare failures✗ Unable to identify problem sources✗ Easy for developers to "game"
Grey-Box <i>(Limited API access to internals developers choose to expose)</i>	<ul style="list-style-type: none">✓ Access predefined outputs only (logits, embeddings)✓ Analyze basic confidence scores / probabilities✓ Limited testing of explicitly exposed features	<ul style="list-style-type: none">✗ Limited insight into internal representations✗ Exact problem mechanisms may remain hidden✗ Unable to modify or fine-tune model
'De facto' White-Box <i>(Broad API access to internals but cannot copy / extract full model parameters)</i>	<ul style="list-style-type: none">✓ Access to internal activations at specific layers✓ Run custom functions on exposed states✓ Limited interventions (exposed activations)✓ Maintain protection of model IP	<ul style="list-style-type: none">✗ Higher latency for multi-step analyses / iterative experiments due to API calls✗ Still restricted to approved intervention points✗ No visibility over weights / attention patterns✗ Limited causal tracing across model components
White-Box <i>(full access incl. weights, activations, gradients, fine-tuning)</i>	<ul style="list-style-type: none">✓ Access model internals loading model locally✓ Trace behaviors to specific components✓ Reveal hidden capabilities via fine-tuning✓ Complete causal analysis of behavior	<ul style="list-style-type: none">✗ No insight into training methods or data✗ Requires model-specific architectural knowledge✗ Demands infrastructure expertise for large models
Outside-the-Box <i>(access to training data and deployment information)</i>	<ul style="list-style-type: none">✓ Comprehensive inspection of model✓ Full intervention at any computational point✓ Evaluate development practices	<ul style="list-style-type: none">✗ Cannot show how model processes information✗ Documentation may be incomplete✗ Requires model access for complete picture



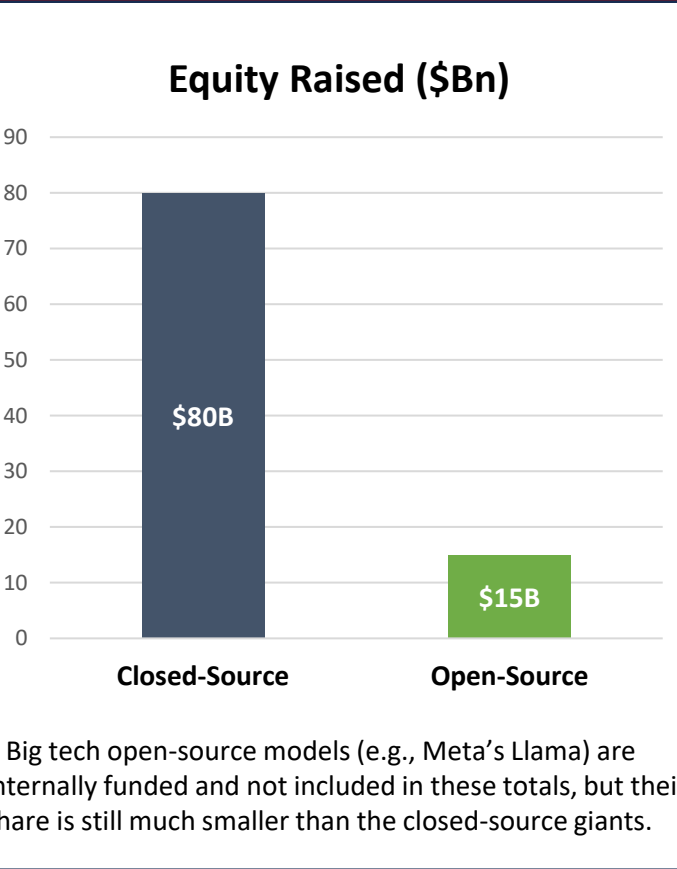
Black-box models dominate, limiting interpretability efforts

Priority to understand most widely used and capable models, but these are predominantly black-boxed

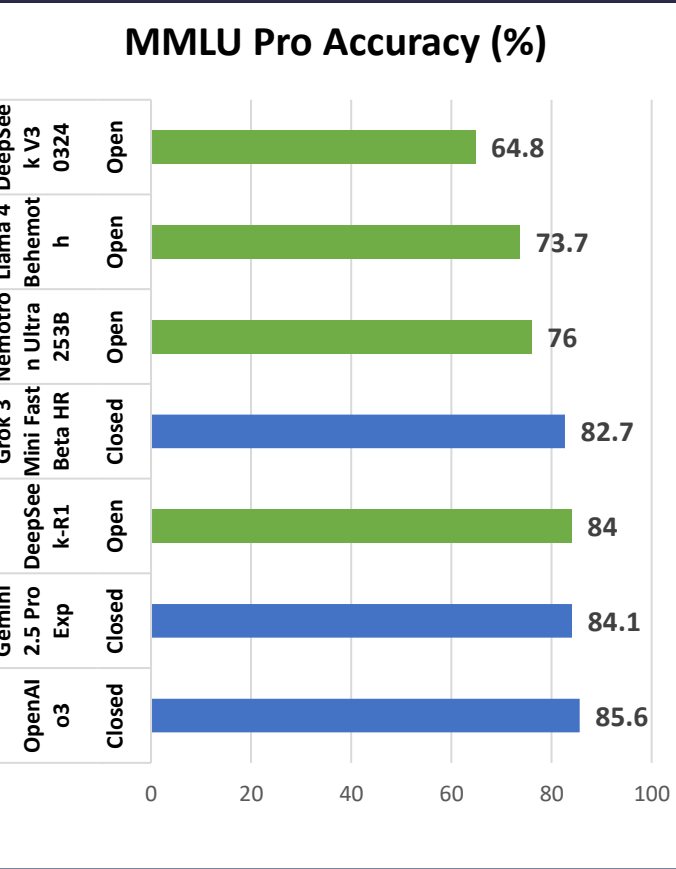
Adoption: Use of Closed-source models predominates with ~90% market share



Investment: Closed-source developers attracted ~5x more funding since 2020



Capability: Closed-source still leads in performance but gap is closing rapidly

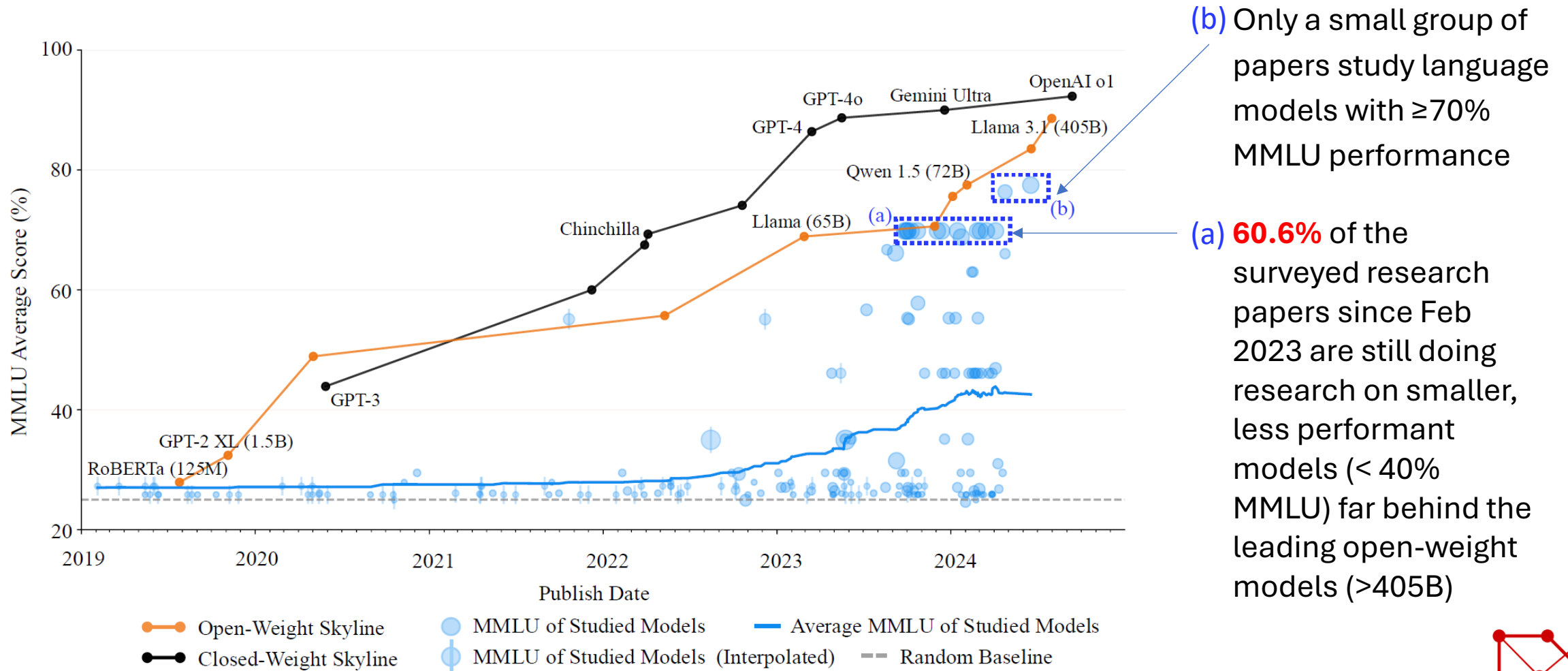


Source(s): CBInsights: The Foundation Model Divide (Jan 2025); GAI Insights: Estimated Market Share of Closed-Source LLM Models (Aug 2024); Vals AI MMLU Pro Benchmark (May 2025); Statista MMLU-Pro Trends (May 2025); Vellum AI Open LLM Leaderboard (Apr 2025)



Even for open models, studies of leading models are lacking

Most interpretability research is done on models that lag far behind SoTA capabilities.
This gap is extended even further when considering leading closed-weight model the performance.



Gap suggests engineering & infrastructural barriers to research

Interpretability research on >70B models is a leap in technical complexity & resource requirements



H100 GPUs	4	➡	~6x	➡	24
Total VRAM	280 GB	➡	~5.7x	➡	1.6 TB
Inference Cost	\$0.80/M tokens	➡	15x	➡	\$12.00/M tokens
Power	~2 KW (4 H100 GPUs)	➡	~8x	➡	~16 KW (24 H100 GPUs)
Network	100 Gbps (single node)	➡	~4x	➡	400 Gbps (multi-node)
Technical Complexity	Single-node tensor Stable at FP4 quantization Most ML tools support single-node		➡		Multi-node pipeline + tensor parallelism Custom kernel optimizations Multi-TB Checkpoint Management Debugging with Distributed logs (takes days) High Error rates with quantization High downtime due to GPU / network failures

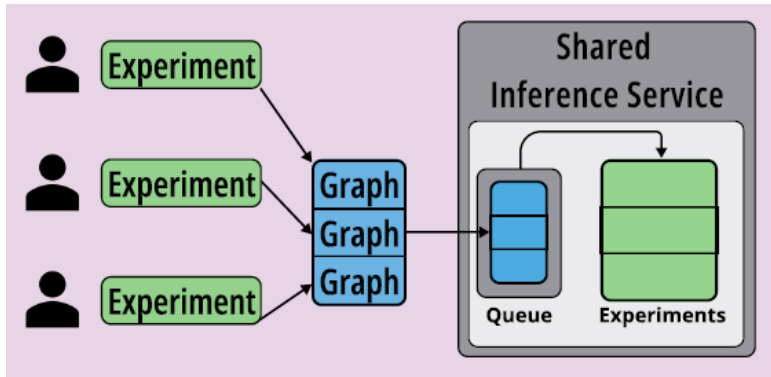
Note: these estimates reflect significant incremental costs involved in mechanistic interpretability experiments (VRAM overhead from activations and gradients, increased token processing, unviable quantization, etc.)



NDIF tackles infra barriers and complexity to study of large NNs

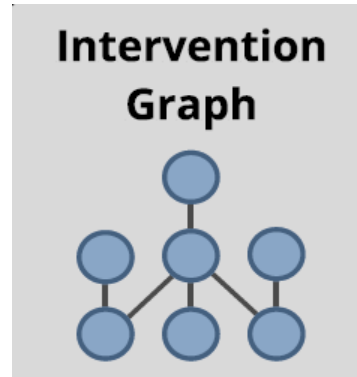
Co-tenanted compute, a new framework to design experiments and a data structure to represent them

NDIF Shared Inference Service minimizes infrastructural barriers



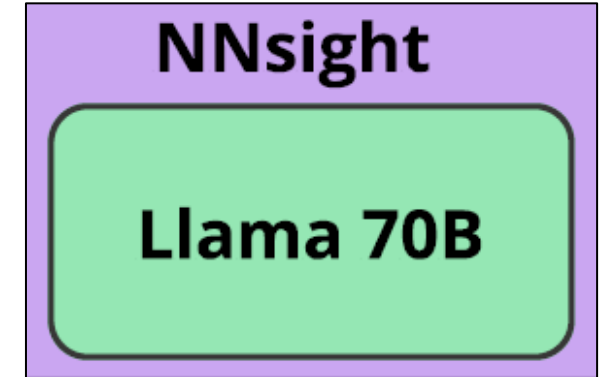
- Multi-user runtime infrastructure designed to ***amortize costs and reduce engineering burden***
- Designed to allow multiple researchers to conduct custom experiments while ***sharing the underlying pre-loaded model instances on shared compute HW***

Intervention Graph structure separates experiment code to run on shared infra



- Data structure holding experiment code, ***dynamically captures data from and runs modifications to the model***. Stored in JSON, and sent to or retrieved from NDIF inference server for execution
- Decouples experimental and engineering code so that the ***complexity of running the model can be tackled separately***

NNsight Framework to define intervention graph and run remotely



- Python library that extends PyTorch to ***define and execution intervention graphs***, introducing deferred remote execution.
- Trace context window provides access to the intermediate inputs and outputs of all modules
- Allows the user to express all the experiments ***as if using PyTorch locally***.



Crucially, NNSight implements and manages intervention graphs at run time between remote & local

- Extends PyTorch with **familiar syntax** for interventions
- Uses **context manager** (with *model.trace()*) for defining experiments
- **No need to manually register/remove hooks** for each experiment
(*Standard PyTorch requires creating custom hooks for each access point*)

- Operations **recorded to intervention graph for later remote execution** within single forward pass
- Creates **proxy objects representing tensor values to be captured** for later interaction
- Marks to **preserve tensor values after execution** with .save()
- FakeTensor Integration enables **experiment debugging without full execution**

- **Access** any model inputs, outputs, and gradients
- Support for **complex interventions** (e.g., activation patching, causal tracing)
- Provides consistent **access to both activation values** (forward pass) **and gradients** (backward pass)
- **Compatible** with a wide range of built-in and custom PyTorch model architectures

- **Same code runs locally** or on NDIF servers by simply adding `remote=True`
- **Debug** on small models locally, then run identical code on larger models to debug effectively
- **No distributed programming needed** - identical API abstracts away model sharding and parallelism



Neural networks can be represented as computational graphs

Computational graph as an alternative representation of NNs to the traditional 'neuron' model

Traditional 'Neuron' Representation blends both data and computation into discrete 'neuron' layers

The traditional neuron concept is a simplified teaching model that combines:

- Parameters (weights & biases)
- Computation (weighted sum & activation)
- Data flow (inputs & outputs)

A "neuron" is described as if it:

- Receives inputs
- Computes a weighted sum
- Applies an activation function
- Sends an output forward

Computational Graph Representation separates these more accurately reflecting how NNs are implemented

There are no "neurons" as unified entities, but tensors and matrix operations (closer to PyTorch):

- **Variable Nodes (Tensors):** Weight matrices by bias vectors by layer; input (from previous layer) & output (to next layer) activation tensors
- **Apply Nodes (Operations):** Matrix multiplication (input \times weights), bias addition, activation function application

Weakly connected, bipartite, acyclic graph with directed edges linking variable and apply nodes:

- Input Tensor (Variable) \rightarrow Matrix Multiply (Apply) \rightarrow Intermediate Tensor (Variable) \rightarrow Activation Function (Apply) \rightarrow Output Tensor (Variable)



Intervention graph is a structured way to modify computational flow

Interventions are modifications to the computational graph that interleave new nodes and edges

An intervention graph modifies a neural network's computational flow by:

1. **Creating a secondary computation graph (C')** with custom operations
2. **Connecting it to the original graph (C)** at specified points
3. **Executing the combined graph** in a single forward pass

The secondary computation graph (C') is comprised of two types of nodes (with corresponding edges):

- **Getters (access hidden states):** Extract tensors from original computation, connecting *variable nodes* in original graph (C) to *apply nodes* in intervention graph
- **Setters (replace hidden states with modified values):** Insert modified tensors back into the original computation, connecting *variable nodes* in intervention graph to *apply nodes* in original graph (C)

Getter: accesses the raw tensor values from specified variable node in layer 16 of the NN

Setter: Modified values are implicitly inserted back (NNsight handles this automatically when a tensor is modified)

```
from nnsight import LanguageModel
model_id = "meta-llama/Meta-Llama-3.1-8B"
lm = LanguageModel(model_id)
```

```
mlp = lm.model.layers[16].mlp.down_proj
neurons = [394, 5490, 8929]
prompt = "The truth is the"
```

Intervention Operation: Sets specified tensor elements controlling "semantic inversion" to a high activation value and save result to 'out'

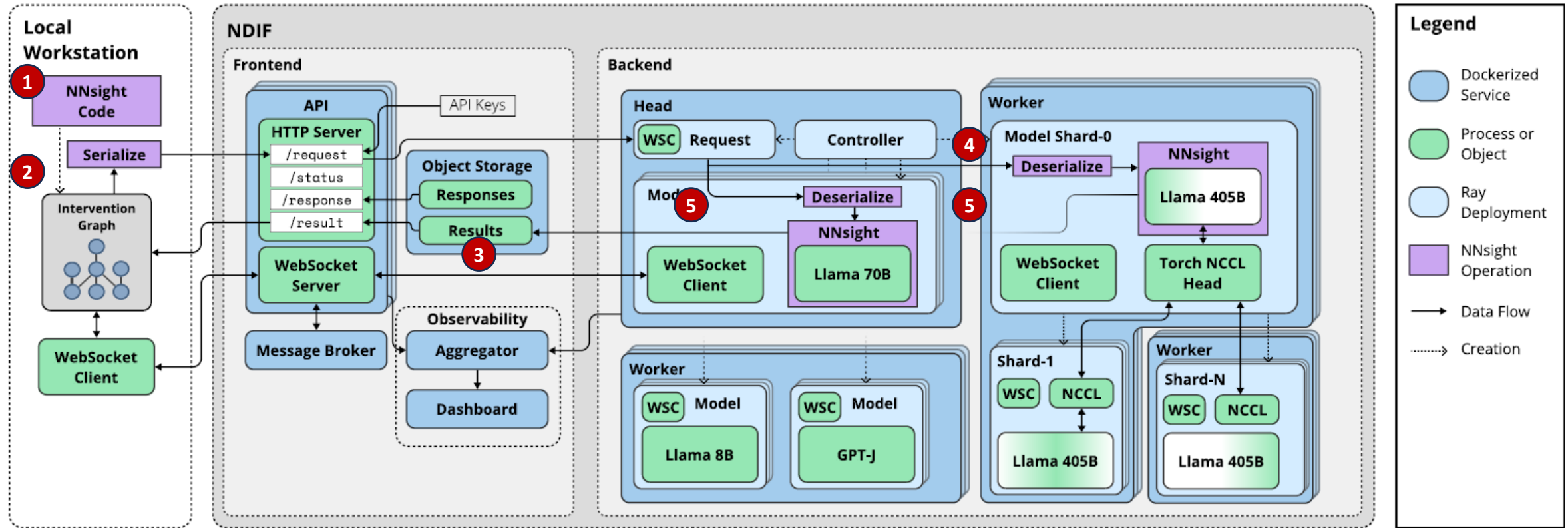
```
with lm.trace(prompt, remote=True):
    mlp.input[:, -1, neurons] = 10
    out = lm.output.save()
```

```
last = out["logits"][:, -1].argmax()
prediction = lm.tokenizer.decode(last)
print(prediction)
```

The most probable token predicted – after the intervention operation – is extracted and decoded to English, resulting in "The truth is the lie"



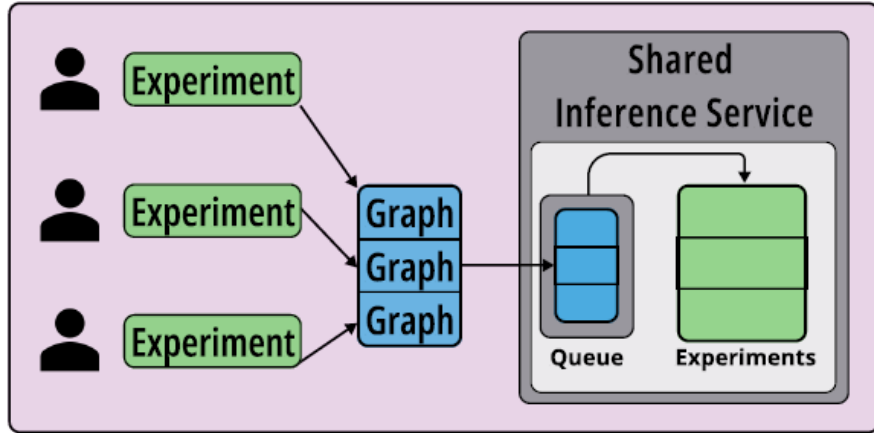
NNsight + NDIF work in tandem to enable interpretability studies



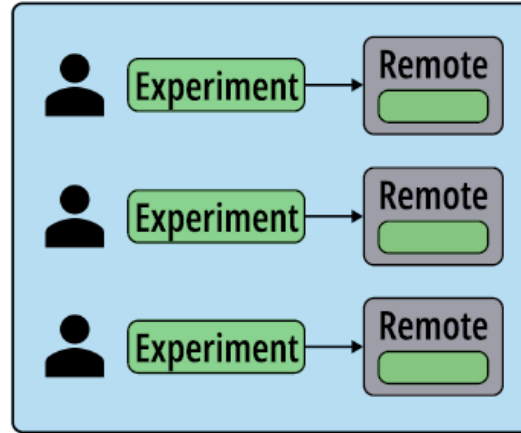
- 1 Researchers write experiment code using the NNsight API
 - 2 Experiment code is converted to an Intervention Graph and sent to the NDIF frontend server
 - 3 The intervention graph is routed to a worker node at NDIF's backend, which hosts multiple models
 - 4 For larger models, (e.g. Llama 3.1 405B), model weights are distributed across shards (tensor parallelism)
 - 5 Results are then gathered and sent to the object store in the NDIF frontend, then pulled to local workstation
- * The graph can either be executed locally, or it can be sent to the remote NDIF servers for execution



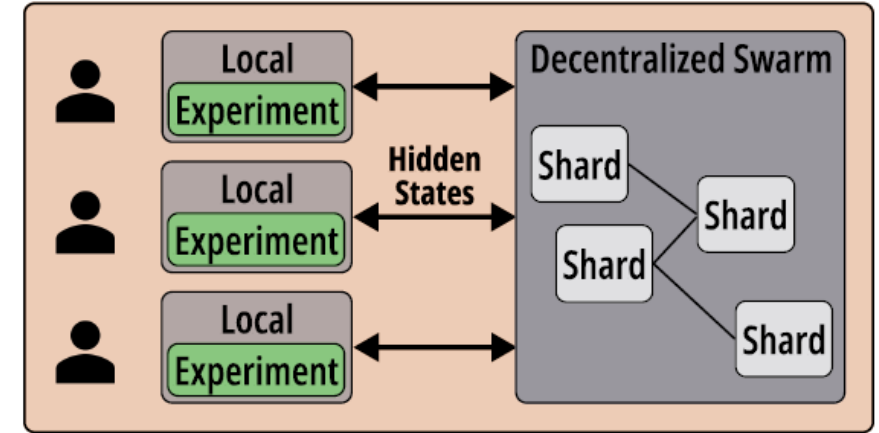
NDIF vs. Traditional Approaches: Eliminating Infra Barriers



NDIF



HPC



PETALS

- | | | |
|--|---|---|
| • Multiple researchers share model instances | • Separate model instance per researcher | • Model distributed across volunteer nodes |
| • Models are pre-loaded, no startup delay | • Extensive model loading times per session | • One-time loading, but network latency |
| • Experiment code separate from model infra | • Researcher manages experiment + infra | • Researcher runs local experiment code |
| • Costs amortized across many users | • Full GPU cluster cost per researcher | • Distributed costs (limited to largest models) |
| • Abstracts model parallelism complexity | • Requires expertise in distributed systems | • Requires P2P networking knowledge |
| • Computations execute entirely on server | • Computations stay on researcher's VM | • Hidden states transfer downloaded locally |
| • Simple API with local/remote compatibility | • Complex distributed programming | • Custom API with limitations on interventions |

NDIF democratizes frontier model research by providing managed shared compute, eliminating the prohibitive costs of traditional HPC and the network bottlenecks of peer-to-peer approaches like Petals



NDIF as Industry Standard: Enabling ‘De Facto White-Box’ Access

At ~90% ‘market share’, it is an imperative for NDIF to seek adoption by SoTA proprietary model providers

“A key limitation is that [NDIF] does not provide access to closed-parameter, proprietary models like GPT-4 or Claude.

Proprietary vendors often (...) wish to maintain secrecy and security of model parameters for business and safety reasons. Our system is designed with these concerns in mind.

We encourage commercial providers to consider adopting our approach as part of their products to enable transparency and facilitate future scientific progress at the frontiers of AI”

Commercial Adoption Pathways

- **Establish NDIF as industry standard** for "de facto white-box" access while protecting proprietary models
- **Implement middleware solutions** that either embed NNSight or connect NDIF protocols to the existing commercial APIs
- **Enable selective transparency** through API extensions exposing specific activations without weight access
- **Bridge the gap** between fully closed systems and open-weight models to support transparency

Multi-Stakeholder Benefits

- **Preserve commercial IP** while demonstrating transparency commitment
- **Enable research** of frontier models powering 90% of deployed AI systems
- **Support regulatory compliance** with emerging AI transparency frameworks
- **Build user trust** through verifiable model inspection and understanding

Strategic Partnership Opportunities

- **Integrate with cloud providers** (AWS, GCP, Azure) as managed service
- **Connect with MLOps platforms / HuggingFace** to streamline adoption
- **Collaborate with hardware manufacturers** (NVIDIA) to optimize performance and embed NDIF/NNSight into model-agnostic toolkits
- **Partner with open model providers** (Meta) to expand ecosystem reach
- **Engage regulators, think tanks** and lobbies to become regulatory standard

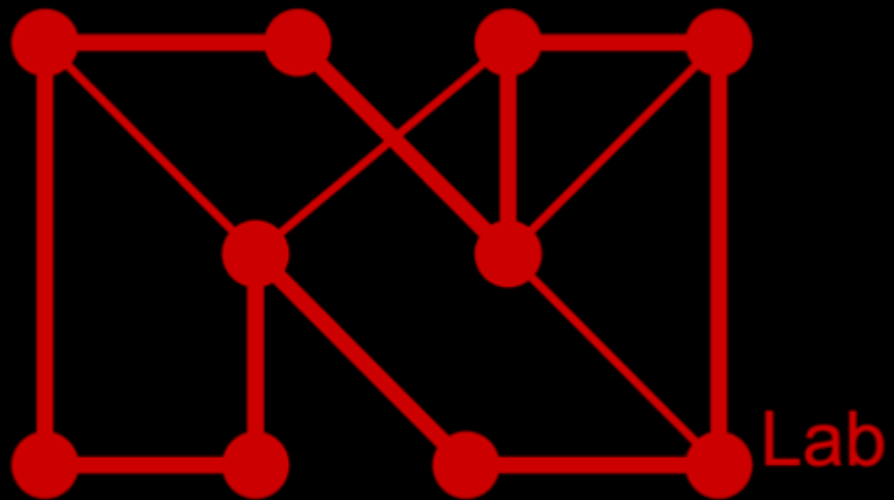


Appendix

Black-box models severely limit interpretability studies

Access Type	Capabilities & Interpretability Goals	Limitations
Black-Box	<ul style="list-style-type: none"> • <i>Test inputs/outputs</i> to assess basic performance • Run <i>manual red-teaming</i> to find obvious flaws • Analyze <i>response patterns</i> to detect biases • Perform <i>statistical analyses</i> to measure fairness • Create <i>simple behavioral tests</i> to map boundaries 	<ul style="list-style-type: none"> • Cannot detect <i>rare or unexpected failure modes</i> • Unable to identify <i>which model parts</i> cause failures • Findings <i>limited by the inputs</i> you think to test • <i>Developers can optimize for test cases</i> without fixing core issues • <i>No way to verify</i> if model <i>explanations</i> are truthful
Grey-Box	<ul style="list-style-type: none"> • Access <i>token probabilities</i> to analyze uncertainty • Examine <i>basic activation patterns</i> to trace influences • Use <i>probability-guided attacks</i> for more efficient testing • Analyze <i>confidence scores</i> to detect potential issues • See <i>limited internal states</i> to improve attribution 	<ul style="list-style-type: none"> • <i>Cannot pinpoint exact neural connections</i> causing problems • Limited ability to <i>uncover hidden capabilities</i> • Unable to use <i>efficient gradient-based</i> testing methods • <i>Cannot modify</i> the model to test how easily safeguards break
De facto White-Box	<ul style="list-style-type: none"> • Run indirect processes via API for <i>structured probing</i> • Use <i>transfer-based attacks</i> to find vulnerabilities systematically • Perform <i>gradient-free optimization</i> for automated testing • Test <i>model manipulation</i> without direct weight access • Run <i>batched experiments</i> to detect internal biases 	<ul style="list-style-type: none"> • Only <i>indirect access through an API</i>, not to actual model code • Testing <i>restricted</i> to what the developer's interface allows • Some <i>analyses take longer</i> due to API limitations • Limitations on certain internal <i>pattern analyses</i> • <i>Harder to trace specific behaviors</i> to exact model components
White-Box	<ul style="list-style-type: none"> • Analyze weights to <i>attribute behaviors to specific components</i> • <i>Fine-tune models</i> to reveal hidden/dormant capabilities • Study <i>internal representations</i> to identify concept formation • Run <i>gradient-based analyses</i> for mechanistic understanding • Perform <i>latent space attacks</i> to test safety boundaries thoroughly 	<ul style="list-style-type: none"> • <i>Requires specialized technical knowledge</i> to analyze effectively • Analyzing large models needs <i>substantial computing resources</i> • No visibility into <i>how the model was trained or on what data</i> • Cannot determine <i>why specific patterns were learned</i>
Outside-the-Box	<ul style="list-style-type: none"> • Review training methodology <i>to trace problems to decisions</i> • <i>Analyze training data</i> to assess bias and quality issues • Study documentation to <i>evaluate development practices</i> • <i>Examine deployment</i> settings to understand real-world risks • Assess developer practices for <i>accountability</i> 	<ul style="list-style-type: none"> • Does not show <i>how the model actually processes</i> information • <i>Documentation may contain gaps</i> or inaccuracies • Hard to connect specific <i>development choices to model behaviors</i> • Companies may still <i>withhold proprietary technical details</i> • <i>Cannot detect emergent properties</i> without examining the model





NORTHEASTERN .