

Charting and Navigating Hugging Face’s Model Atlas

Eliahu Horwitz

Nitzan Kurer

Jonathan Kahana

Liel Amar

Yedid Hoshen

School of Computer Science and Engineering
The Hebrew University of Jerusalem, Israel

<https://horwitz.ai/model-atlas>

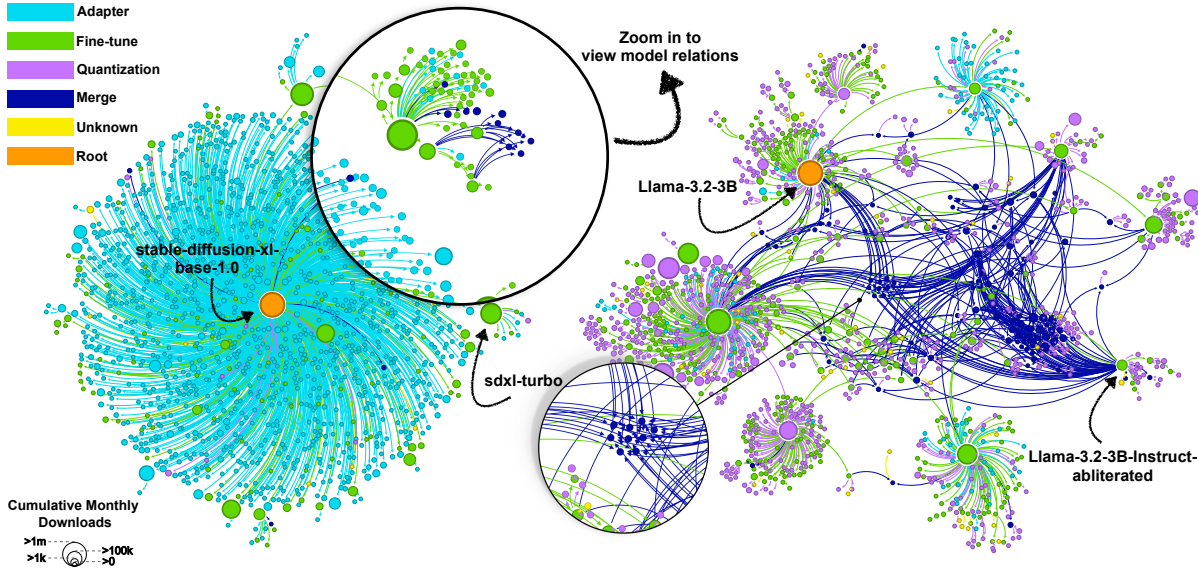


Figure 1. *The model atlas - Stable Diffusion vs. Llama*: The model atlas visualizes models as nodes in a graph, with directed edges indicating transformations (e.g., fine-tuning). This figure shows the top 30% most downloaded models in the Stable Diffusion and Llama regions. Node size reflects cumulative monthly downloads, and color denotes the transformation type relative to the parent model. Please zoom in to see the detailed model trajectories. We observe that the Llama region has more complex structure and a wider diversity of transformation techniques (e.g., quantization, merging) compared to Stable Diffusion. Note that node position is optimized for clarity [20] and does not directly reflect distance between model weights. *Zoom in to view edges, best viewed in color.*

Abstract

As there are now millions of publicly available neural networks, searching and analyzing large model repositories becomes increasingly important. Navigating so many models requires an atlas, but as most models are poorly documented charting such an atlas is challenging. To explore the hidden potential of model repositories, we chart a preliminary atlas representing the documented fraction of Hugging Face. It provides stunning visualizations of the model landscape and evolution. We demonstrate several applications of this atlas including predicting model attributes (e.g., accuracy), and analyzing trends in computer vision models. However, as the current atlas remains incomplete, we propose a method for charting undocumented regions. Specifically, we identify high-confidence structural priors based on dominant real-world model training practices. Leveraging these priors, our approach enables accurate mapping of previously undocumented areas of the atlas. We publicly release our datasets, code, and interactive atlas.

1. Introduction

Navigating millions of publicly available neural networks is a challenging task that could greatly benefit from an atlas: a structured representation capturing how models evolve, what tasks they solve, and how well they perform. Such an atlas would allow users to easily discover and reuse existing models rather than training new ones from scratch, saving resources and reducing environmental impact. Moreover, it could serve as a snapshot of the entire machine learning landscape, facilitating comparisons across techniques and modalities and highlighting emerging trends. In this paper, we propose a method to chart this atlas and demonstrate its practical utility for navigating large model collections.

To establish the untapped potential of atlases of large model repositories, we first chart the documented regions of Hugging Face (HF). Its visualizations reveal the intriguing and complex structures of the machine learning model landscape. We use these visualizations to analyze recent trends in computer vision models and compare them across different modalities. For example, Fig. 1 showcases a portion

of our atlas, highlighting trend differences between Stable Diffusion [34] and Llama [9]. It is clear that Llama-based models employ more diverse training practices and dynamics, including quantization and model merging, resulting in more complex structures. We further demonstrate how the atlas structure can predict both model tasks and their accuracy. For example, we improve the prediction of the TruthfulnessQA metric of models derived from Mistral-7B by simply observing their nearest neighbors in the atlas.

In reality, model documentation only provides an initial glimpse of the atlas, while the full picture remains elusive. Existing model metadata (e.g., Model Cards [28] or configuration files) are often incomplete and lack critical details about model task [16, 22], accuracy (Fig. 3), and origin [18]. We tested the completeness of HF documentation and found that approximately 60% of the 1.5 million models lack any. This is not only confined to niche models, e.g., MobileNet, has over 70 million monthly downloads yet only 5 documented fine-tunes. Consequently, prior research has explored model tree recovery, representing models as nodes and their relation (e.g., fine-tuning) as edges. Most methods use probing [50] or direct weight inspection [18, 51, 52], but they typically rely on unrealistic assumptions (Sec. 4) and fail to scale to the real-world.

To this end, we develop an atlas charting method for real-world repositories. We identify common real-DAG patterns neglected by current approaches. In particular, we analyze the harmful effects of duplicate or near-duplicate models on existing methods. This leads us to identify and use priors on model relations that allow us to predict the correct edges. Also, the increasing popularity of model merging violates the traditional tree-structure assumption. We thus represent the atlas as a non-tree directed acyclic graph (DAG). Our proposed algorithm, informed by these assumptions, effectively recovers the atlas structure of HF, achieving substantial performance improvements over baseline methods.

Our primary contributions include:

1. Demonstrating that large-scale model repository atlases are more interact than previously suggested.
2. Showcasing practical applications of the atlas, including model trend analysis and attribute prediction.
3. Introducing an effective method for atlas charting under realistic, in-the-wild conditions by identifying and leveraging structural patterns in real-world DAGs.

2. Related works

Weight-space learning. The emerging field of studying models based on their weights is called Weight-Space Learning [10, 13, 36–38, 42]. Research in this area generally falls into 3 categories: i) *Learning weight representations* to recover the functionality and performance of models [15, 16, 21–23, 25, 29, 30, 33, 35, 53]. ii) *Developing generative models for weights*, enabling the synthesis of new

model parameters [2, 8, 11, 14, 32, 39, 47]. And iii) *Recovering weight trajectories*, which involves tracing the evolution of model parameters over time [5, 6, 17, 18, 45, 49–52].

Analyzing model repos. As large model repos have only grown recently, studying them is a new field. Some works focused on analyzing the hosted datasets [48] and examining broader development trends [12, 31]. Differently from these papers, we analyze the heritage structure of models repositories and show some motivational applications.

Model heritage recovery. A recent line of work [18, 50–52] suggested a way to recover the model heritage structure from model weights. They were able to recover model lineage from curated or synthetically generated model populations. In this work, we improve on these ideas by extending them to real-world model repos. This requires rethinking the assumptions made by heritage methods.

3. Analyzing Hugging Face’s model atlas

The atlas \mathcal{M} consists of a set of models. Each model $m \in \mathcal{M}$ has weights (parameters) w . It also has attributes including: the model’s name, its upload time (t), parents (P) if known, training hyperparameters, performance on specific benchmarks, or download count. A child model m_{ch} is formed by transforming the weights of the parent model m_{pa} through operations like finetuning, quantization, pruning, or merging [3, 40, 41, 47]. We denote this by adding a directed edge (m_{pa}, m_{ch}) to the atlas. A model with no parents is called a *source* model. As the weight transformation is causal, i.e., a model cannot be a descendant of itself, the atlas structure is a directed acyclic graph (DAG). Differently from previous work [18, 52], we do not assume the DAG is a tree as some models can have multiple parents, particularly *merges*. We call each connected component in the DAG an atlas region.

In this section, we analyze Hugging Face, the largest public model repository. It currently hosts over 1.5 million models, with over 100,000 new models added each month, totaling over 17 PB [44] of storage. While HF encourages model documentation through Model Cards [28], the vast majority are poorly documented and over 50% of models have no model card at all. In particular, HF has labels for only 400,000 directed edges. Using these labels we construct the initial atlas and demonstrate its promise for analyzing recent trends, measuring model impact, predicting model accuracy and other attributes. We make this atlas publicly available via an interactive demo website.

3.1. The atlas reveals trends in machine learning

In Fig. 2 we present a part of the HF atlas. It contains over 63,000 models spanning 28 CCs with a total of over 65,000 edges. We use this powerful visualization to compare machine learning trends in computer vision (CV) and natural

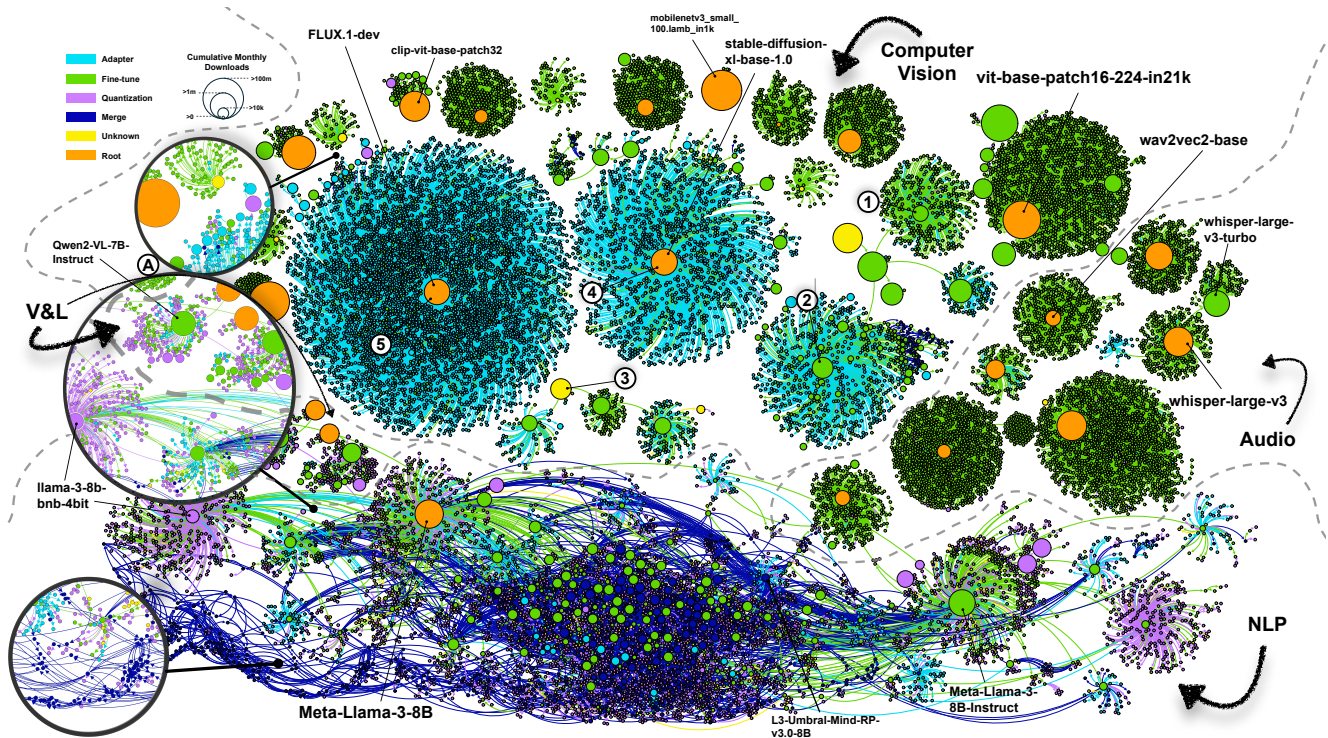


Figure 2. **The Hugging Face atlas:** While this is a small subset (63,000 models) of the documented regions of HF, it already reveals significant trends. *Depth and structure.* The LLM connected component (CC) is deep and complex. It includes almost a third of all models. In contrast, while Flux is also substantial, its structure is much simpler and more uniform. *Quantization.* Zoom-in (A) highlights quantization practices across vision, language, and vision-language (V&L) models. Vision models barely use quantization, despite Flux containing more parameters (12B) than Llama (8B). Conversely, quantization is commonplace in LLMs, constituting a large proportion of models. VLMs demonstrate a balance between these extremes. *Adapter and fine-tuning strategies.* A notable distinction exists between discriminative (top) and generative (bottom) vision models. Discriminative models primarily employ fine-tuning, while generative models have widely adopted adapters like LoRA. The evolution of adapter adoption over time is evident: Stable-Diffusion 1.4 (SD) (1) mostly used full fine-tuning, while SD 1.5 (2), SD 2 (3), SD XL (4), and Flux (5) progressively use more adapters. Interestingly, the atlas reveals that audio models rarely use adapters, suggesting gaps in cross-community knowledge transfer. This inter-community variation is particularly evident in *model merging*. LLMs have embraced model merging, with merged models frequently exceeding the popularity of their parents. This raises interesting questions about the limited role of merging in vision models. For enhanced visualization, we display the top 30% most downloaded models. *Zoom in to view edges, best viewed in color.*

language processing (NLP) models. See App. B for additional visualizations.

Depth. Fig. 3 shows that NLP models have a wide range of DAG depths, while in CV, nearly all models have a direct edge to the root foundation model. This suggests that the CV community puts more focus on new foundation models, while the NLP community often embraces iterative refinement. As an example, in Fig. 2 The LLM CC exhibits significant depth and complexity, representing almost a third of the models. In contrast, while Flux is also substantial, its structure is much simpler and more uniform.

Quantization. Fig. 2 shows that the vision community has not adopted the use of quantization (less than 0.15% of all models in this pool). Meaning, vision models have not yet reached the scale where inference time is too costly and

quantization is essential. However, using the atlas we can spot an emerging change of trend. Flux [24], one of the newest and largest generative models, is the rare exception in terms of quantizations among vision models. This indicates that image generation models may have just reached the scale where quantization is valuable.

Fine-tuning vs. adapters. A striking difference between generative and discriminative models is that the vast majority of generative models use adapters (e.g., LoRA) while almost all discriminative models use full fine-tuning (see Fig. 2). The atlas also allows us to track how trends evolve. In Fig. 2 we can trace adapters over time and see they are becoming more prevalent. For example, only 50% of older models (such as SD1.4) use parameter-efficient adapters, while in newer versions it is much more common. The

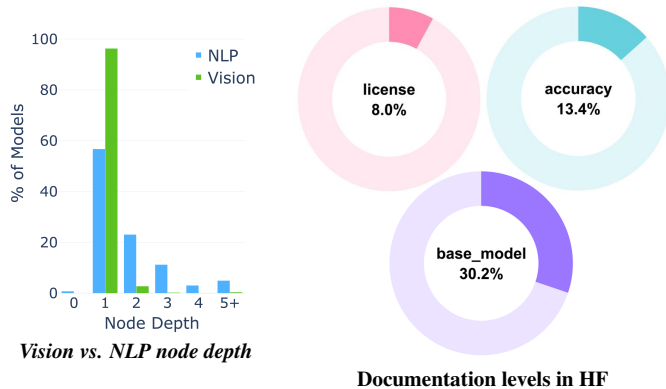


Figure 3. **Left:** By analyzing over 314k models, we found that over 96% of CV models are situated one node away from the root, while only 55% of NLP models have this shallow depth. Over 5% of NLP models have depth of at least five nodes. This shows that NLP models are much deeper than CV models, suggesting the NLP community embraces iterative refinement over new moving to the latest foundation models. **Right:** A significant portion of models on Hugging Face suffer from poor documentation quality.

newest generative models, such as Flux or Llama 3 [9], have much more parameter-efficient trained models.

Merging. The rate of model merging [40, 41, 43, 46, 47] in NLP is approximately 35 times higher than that of vision models (see Fig. 2). Merged NLP models had, on average, 30% higher influence (total descendant downloads, see Sec. 3.2) than their non-merged sibling models. While not a causal relationship, this hints that more research into merging vision models can be fruitful.

3.2. Model attribute prediction using the atlas

Local atlas regions contain related models. Therefore, the atlas may also be useful for predicting missing model attributes, including task, accuracy, license, missing weights, and popularity.

Model accuracy prediction. Currently, fewer than 15% of HF models have accuracy details (see Fig. 3). This is a serious limitation as users typically want to select accurate models for their task. Even worse, there are very few models that report the same accuracy metrics on the same datasets. Surprisingly, the challenge is not limited to unpopular models, as many popular models do not report accuracy metrics in easily processed form. Here, we demonstrate that the model DAG can help predict model accuracy. To test this, we used the Mistral-7B CC which contains 17.5k models. Only 300 models were labeled with their performance on the *TruthfulQA (0 shot) (mc2)* metric (see App. D for data collection details). We predict the accuracy of each unlabeled node using the average of its K nearest neighbor nodes, where distance is measured by path length on the undirected version of the model DAG. The results in Tab. 1a

Table 1. **Atlas-based documentation imputation:** Using atlas structure improves prediction of model accuracy and other attributes, compared to naively using the majority label. In (b), we report the prediction accuracy.

		(a)			
		Method	MSE	MAE	Correlation
TruthfulQA (0-shot, mc2)	Baseline	100.217	8.541	-	
	Ours 1-NN	32.830	3.247	0.856	
	Ours 2-NN	28.720	3.235	0.864	
	Ours 3-NN	25.544	3.147	0.877	
	Ours 5-NN	23.512	3.093	0.885	
Hellaswag (0-shot)	Baseline	95.000	7.500	-	
	Ours 1-NN	30.000	3.000	0.860	
	Ours 2-NN	27.000	2.900	0.870	
	Ours 3-NN	24.000	2.800	0.880	
	Ours 5-NN	22.000	2.700	0.890	
		(b)			
Attribute	Graph Avg.	Hub Avg.			
pipeline_tag	0.60	0.79	(+0.19)		
library_name	0.81	0.84	(+0.02)		
model_type	0.66	0.81	(+0.15)		
license	0.49	0.85	(+0.35)		
relation_type	0.61	0.80	(+0.19)		

show that the neighbors of a model can predict its accuracy.

Predicting other attributes. Atlas structure can predict other model attributes besides the accuracy. Here, we use model *hubs*, which are sets of sibling leaf models (79% of models are members of some hub). The idea is that models in hubs are very similar to one another, and therefore we can complete their missing values by the majority class of labeled nodes within the same hub. We tested hub-based predictions on 5 attributes that are often missing on HF, including: license, model and inheritance types (see Fig. 3). Tab. 1b compares this method with the graph-level majority label. It improves significantly, including a 35% gain in accuracy for license prediction and 19% improvements for both inheritance type and pipeline tag prediction.

Better model impact measurement. There are several ways of measuring the impact of a model. For example Hugging Face uses likes, trends (based on an undisclosed proprietary algorithm) and downloads. These metrics are somewhat myopic, as they measure the direct popularity of models but not the popularity of their descendant models. In fact, our graph analysis reveals that for 50% of non-leaf nodes, the total downloads of their descendants exceed their own individual downloads. This is partially due to the popularity of quantized child models and to incremental improvement of child models, e.g., finetuning or merging. Our analysis further showed that for non-leaf models the sum of descendant nodes downloads exceeds those of the model it-

self in most cases (often by large margins). This suggests that simple model download counts underestimate the influence of the parent model.

We can therefore use the atlas to introduce a new model impact metric: `sub_tree_downloads`. We calculate this metric by summing the downloads of the model node and those of all its descendants. This number describes how many downloads this model causally affected. It has important applications to intellectual property rights, as the models and data used to train the target models affected all of its downstream downloads. It also quantifies the social impact of the biases of this model.

Restoring removed models. There are legal and commercial reasons for removing models from repos. Some models have been removed from the repository over time, affecting the integrity of the model DAGs. Out of 33,870 identified source models, 1,612 are missing due to deletion. A notable case is `runwayml/stable-diffusion-v1-5`, which has 3,038 broken references. To preserve the integrity of the dependency structure, the missing node was reconstructed and its connections manually restored. Using methods such as Horwitz et al. [17] also allows restoration of the weights.

4. Charting the atlas structure

While the preceding section highlighted the importance of the model atlas, in practice, over 60% of it is unknown. This is primarily because model uploaders frequently do not provide parent model information, as illustrated in Fig. 3. Atlas charting aims to recover the missing edges. Recent methods [18, 50–52] have tackled some aspects of this task, but were not applied to fully in-the-wild settings. As such, some of their key assumptions (e.g., tree-like structures) do not hold in real-world repos. Additionally, some methods [50] require running each model on multiple inputs, which is computationally infeasible for millions of models. Here, we introduce an approach suitable for in-the-wild atlas charting.

Charting typically begins by measuring pairwise model distances, which requires a representation for each model as well as a distance function. Common model representations include: weights [18, 50–52], activations [50], gradients [51], outputs [22, 26], metadata [27] or a combination of the above [51]. We represent each model i by its weights w_i . We measure the distance between two models using the Euclidean distance of their representations:

$$D_{ij} = \|w_i - w_j\|^2 \quad (1)$$

The core challenge is charting the atlas structure from this distance matrix. While previous approaches used simple, generic solutions, such as: greedy nearest neighbor or minimal directed spanning tree, this does not generalize to in-the-wild model atlases. The reasons for this include: i) models with the nearest weight distance are not always connected by an edge ii) the model atlas structure is generally

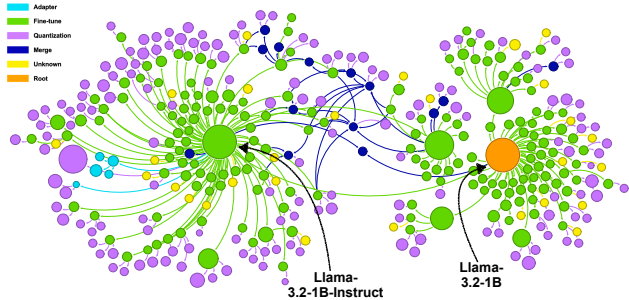


Figure 4. **Quantizations are leaves:** Our analysis of over 400,000 documented model relationships reveals that 99.41% of quantized models are leaf nodes. This figure shows this for a subset of the Llama-based models. Indeed, quantized models (magenta) are nearly always leaf nodes, corroborating the statistical finding.

not a tree. Instead, we propose a set of charting rules motivated by special patterns in model repos: duplication, quantization near-duplications, checkpoint trajectories, hyperparameter search, and merging. Identifying and handling these patterns can dramatically improve charting accuracy.

Naive strategy. Our preliminary strategy first splits models into non-overlapping subsets using a standard clustering algorithm on the weight distance matrix. Previous works [13, 16, 18] showed that each subset corresponds to the nodes within a single connected component of the atlas DAG, i.e., we can recover the structure of each connected component separately. The naive strategy first assigns source nodes (this will be elaborated on below), then, at each step, it looks for the unassigned node with the shortest distance to the assigned nodes, and connects the two with a directed edge. The algorithm stops when there are no more unassigned nodes.

Duplicates and near-duplicates. Real-world model repositories contain many duplicates and near-duplicates. Exact duplicates occur when users download a popular base model and re-upload it without modification. In this, both the original and duplicates have *identical* distances to all models. This makes the greedy distance-based algorithm to arbitrarily decide between the true parent and its duplicates, which obviously reduces accuracy. To mitigate this, we identify models with zero-distance as exact duplicates, and retain only a single representative instance. We achieve this by designating duplicate nodes as leaves.

Model quantization compresses models for reducing memory and storage as well as accelerating their inference. It transforms a model into a near-duplicate, as the values of each weight are typically very similar yet not identical to the original one. This creates a similar ambiguity to exact duplicate models, as the distance of the original and quantized model to any other model are almost the same. Moreover, many models have multiple quantizations, which greatly increases this ambiguity.

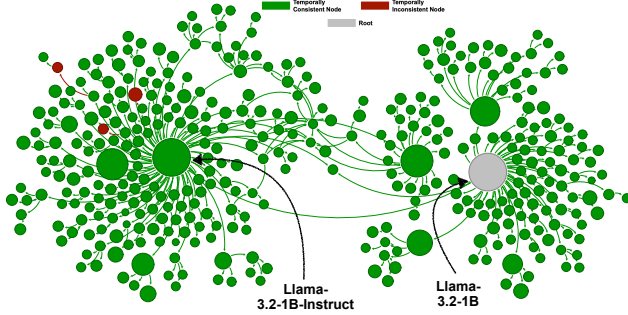


Figure 5. *Temporal dynamics indicate edge directionality*: We analyzed over 400,000 documented model relationships and observed that in 99.73% of cases, earlier upload times correlate with topologically higher positions in the DAG. Here, we visualize this trend on a subset of the Llama model family. Green nodes indicate models where earlier upload times align with topological order, while red nodes represent exceptions to this trend. The source (in gray) vacuously satisfied this assumption. It is clear that nearly all nodes satisfy our assumption.

To overcome this ambiguity, we identify a pattern of quantization models from real-world atlases. Our analysis reveals that 99.41% of quantized models on HF are leaf nodes, i.e., they have no child models. See Fig. 4 for an illustration on the Llama 3.2-1B CC. Intuitively, unquantized models have better performance than their quantized versions, and practitioners typically use the highest-performing models for further fine-tuning or merging. Fortunately, detecting quantization is often straightforward. Quantized models have reduced precision data types (e.g., int8, float16 instead of float32) and fewer unique weight values. Therefore, we address quantization by identifying quantized models and designating them as leaves.

Duplication Assumption

Exact duplicates and quantized versions of models are treated as leaf nodes in the atlas.

Temporal dynamics for inferring fine structure. While the weight distance between models is inversely correlated to the likelihood of having an edge between them, it does not reveal the direction of the edge. Previous approaches predicted edge direction using supervised learning methods or by unlearned metrics, such as, gradients or weight kurtosis. However, we find that these methods often do not generalize effectively to real-world model repositories. A key advantage of real-world repositories, unused by prior works, is the availability of model creation or upload timestamps. We find that the vast majority of parent models have earlier timestamps than their children, specifically, in the HF model repository this occurs for 99.73% of all observed parent-child pairs. We visualize this phenomenon for the

Llama-3.2-1B CC in Fig. 5. This is a powerful temporal constraint on atlas structure, and in particular identifies the source nodes as the earliest ones.

Another limitation of the weight distance is that it may have limited sensitivity in fine-grained atlas structures. Two common scenarios that challenge charting methods based on weight distance are hyperparameter sweeps and checkpoint trajectories. In hyperparameter sweeps, multiple models are trained from a single parent, each with different hyperparameters. Minor hyperparameter changes can result in sibling models having nearer weight distances to each other than to their common parent, defying the nearest neighbor algorithm. Conversely, checkpoint trajectories represent a sequence of models saved during a single training run. In this case, edges should connect consecutive checkpoints in a chain-like manner, rather than connecting each checkpoint directly to the initial model (as in the hyperparameter sweep scenario). We term the former pattern a “fan” pattern and the latter a “snake” pattern, illustrated in Fig. 6, and find that the weight distance often confuses these two patterns.

Our key observation is that the temporal model weight evolution can discriminate between these patterns. In snake patterns, temporal proximity strongly correlates with weight proximity, due to the sequential evolution. However this is not the case in fan patterns, where the closest siblings are not necessarily the closest in time. This leads to a simple yet effective decision rule. If the weight distance of a model to its K nearest neighbors is highly correlated to their absolute temporal distance - we classify the pattern as a snake, otherwise a fan. Then, we connect the node to one of its nearest neighbors depending on the detected structure. In snakes, we connect the node to its first nearest neighbor based on weight distance, while in fans we connect the node to its oldest nearest neighbor, aiming for the fan’s origin.

Temporal Dynamics Assumption

Models with earlier timestamps are assumed to be topologically higher in the atlas. For snake patterns, the parent is the closest preceding model; for fan patterns, it is the earliest model among the top K weight-based neighbors.

Merged models. Prior work assumes that models lie in a directed tree, however, as merged models have multiple parents, they have more than one incoming edge. This cannot be described by trees, but can be described by non-tree DAGs. We found that many merged models are created using a few popular libraries which typically document the parent nodes in the created model card, hence the multiple incoming edges are often known. The challenge left is to chart a DAG instead of a tree. We propose a greedy charting algorithm. We first sort all nodes by upload times. Then,

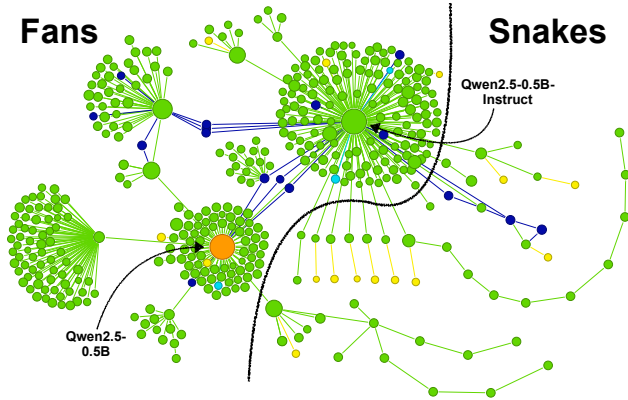


Figure 6. *Snake vs. Fan patterns*: Snake patterns often arise from sequential training checkpoints, while fan patterns typically result from hyperparameter sweeps. In both structures the model weight variance is low. However, in snake patterns the weight distance has high correlation with model upload time, whereas in fan patterns the correlation is lower. Note colors are the same as Fig. 4

we process each node in temporal order, predict its parents, and add those edges to the graph. In the case where the parents are known, we use these edges instead of predicting them. The algorithm stops where there are no more nodes to process. The runtime complexity is $O(n^2)$, where n is the number of models. It is much faster than [18] (which is typically $O(n^3)$). In practice, we recover graphs with thousands of nodes in seconds (see App. A).

Structural Assumption

The atlas structure is a non-tree directed acyclic graph (DAG).

The most expensive part of the algorithm in terms of compute and storage is the distance matrix calculation, as it scales with the number of network weights. Therefore, we propose to represent each model by subsampling the full weight representation to a much smaller number of neurons. Specifically, we find that keeping just 100 is typically sufficient for significantly speeding up the runtime while keeping accuracy nearly unchanged, see App. A for details.

Full algorithm. We present the full algorithm in python-like pseudo code in Alg. 1. The input is a set of models M , each with weights w , creation time t , known parents P (or $P = None$), and whether it is quantized q . The algorithm has 3 hyperparameters: K - the number of neighbors, K_{th}, ρ_{th} - thresholds for snake-fan classification. The output of the algorithm is the predicted parents for each node $m.P$. Note that we can run the algorithm online with a small modification, by computing the distance function and kNN as new models arrive. See App. A for implementation.

Algorithm 1 Atlas structure recovery (psuedo-code)

```
# Input: M = {(w, q, t, P)}, K, K_th, rho_th
# w = weights, t = creation time, P = known parents
M = M.sort(key=m.t) # Sort M by timestamp t
D = compute_distance_matrix(M)
D.diag = inf

for i, m in enumerate(M):
    if is_quantized(m.w): # Set quantized model as leaf
        node
        D[i, :] = inf
    if m.P: # Use existing parents if available
        continue
    for j in range(i+1, len(M)): # Deduplication
        if D[i, j] == 0:
            D[j, :] = inf
            M[j].P = M[i].P

# Find k-NN
k_ind = argsort(D[:, i])[1:K]
temporal_k_ind = argsort(D[i+1:, i])[1:K]
k_dist = D[k_ind, i]
k_times = [M[k].t for k in k_ind]
temporal_k_times = [M[k].t for k in temporal_k_ind]

# Check spread of distances
if k_dist[-1] - k_dist[0] > K_th:
    m.P = temporal_k_ind[0]
else:
    corr = correlation(k_dist, abs(k_times - m.t))
    if corr > rho_th: # Snake
        m.P = temporal_k_ind[0]
    else: # Fan
        m.P = temporal_k_ind[argmin(temporal_k_times)]
```

4.1. Charting results

Datasets. We created an atlas charting dataset based on the ground truth model relation found on the “hub-stats”¹ dataset. It consists of 3 atlas connected components: Qwen2.5-0.5B (Qwen), Llama-3.2-1B (Llama), and Stable-Diffusion-2 (SD). We also created another attribute prediction dataset for the experiments in Sec. 3. Its task is to predict 6 attributes: accuracy, pipeline_tag, library_name, model_type, license, relation_type. The ground truth was obtained by a combination of model metadata and running an LLM on the model cards. We detail the dataset creation process in App. D and will publicly release the datasets.

Baselines. We compared a weight-agnostic classical method and a newer method that depend on weights. *MoTher* [18] - a recent weight-dependent method. Importantly, it predicts edge direction using weight kurtosis and relies on the structure being a tree. We could not scale other weight-based methods to operate on our in-the-wild dataset. We also included the following structure-only baselines: *random* edge assignment, *majority vote* (assigning all nodes to the parent with the highest out-degree), and *Price’s model* [7]. The latter is a preferential attachment algorithm [1], popular for citation networks, that assigns edge probabilities based on node out-degree. It is essentially a probabilistic version of majority vote that allows for new “hub” formation. Note that all baselines (except random), assume

¹huggingface.co/datasets/cfahlgren1/hub-stats

Table 2. *Atlas recovery results*: Our method outperforms the baselines by a significant margin, even for in-the-wild models.

Method	Qwen	Llama	SD
Random - root	1.77	1.84	0.22
Random	0.98	0.67	0.75
Majority	15.03	25.00	36.75
Price’s [7]	2.28	5.08	8.50
MoTher [18]	32.81	19.32	50.51
Ours	78.87	80.44	85.10

the structure is a tree, not a DAG. To accommodate the majority and Price’s models, we assumed a graph stem with known edges, containing 10% of the nodes in the connected component. We evaluated all methods by their accuracy on the remaining 90% of models.

Baseline comparison. We use our method to chart three atlas connected components. As we can see in Tab. 2, the baselines perform better than random but still poorly. The majority approach models the graph as depth-1 centered around the root. This obtains low but non-trivial results, as this model is too simplistic. Price’s model is more powerful and can create a more realistic looking structure. However, it fails as it is data agnostic and connects the wrong nodes. MoTher performs better on some DAGs, but its directional prediction is weaker and it lacks our other priors, resulting in lower performance. In the App. E we visualize the predictions and breakdown of our methods’ errors.

Ablation. Tab. 3 demonstrates that all components of our method have value. The greedy method improves over the exact spanning tree algorithm as it allows us to handle non-tree DAGs which include model merges, and always find a valid solution. For example, the Edmonds’ algorithm failed to return a result on the SD split as it did not find a valid tree. Our quantization prior has a major effect for DAGs with many quantization near-duplicates, such as Llama (+44%). The fans and snakes prior improves all DAGs, especially in hub-like DAGs such as Stable-Diffusion (+13%). In the App. A we ablate approximating the full weight distance with a subset of neurons. Indeed, 100 neurons results in minimal accuracy loss while being 2 magnitudes cheaper.

5. Discussion and limitations

Merged models. Our atlas can handle DAG structures, including models with multiple parents (merged models). Model documentation often specifies if a model is a merge and identifies its parents. For a model that does not have this information, our method is unable to identify if it is a merge or predict its parents from its weights alone. Addressing this remains a topic for future research.

Distilled models. Our atlas currently does not represent distillation relationships between models. Capturing these

Table 3. *Subtractive ablation*: We remove each of our assumptions individually, indeed, we see that each one contributes to our final high recovery accuracy.

	Method	Qwen	Llama	SD
Ours	– Greedy Alg.	77.34	78.98	Failed
	– Quantization	70.57	36.59	84.85
	– Deduplication	67.85	75.28	88.76
	– Temporal Consistency	75.47	80.13	80.86
	– Fans vs. Snakes	75.09	76.03	71.72
	– Merges	76.95	76.61	84.85
	Ours	78.87	80.44	85.10

relationships could be valuable. Our preliminary investigations suggest that analyzing the differences between parent and child models can reveal training dataset and model properties, consistent with findings in [13, 16, 19]. Initial experiments indicate this also applies to individual neurons, suggesting potential for representing student-teacher relationships in distilled models. We leave this to future work.

Additional visualizations and analysis. Beyond visualizing model relations, we can use the atlas to visualize other attributes and dynamics of model repositories (see App. B).

Intellectual property (IP) tracking. Another important application of the full model atlas is intellectual property tracking. Once a model has a particular license or uses data with particular right, all its descendants may be affected. Completing the atlas can therefore improve IP integrity, and prevent unauthorized use.

Other model repositories. While this work focuses on HF due to its popularity and diversity, it is relevant to other model repositories, both public and private. A straightforward example is a machine learning researcher’s file system, or more broadly all models trained by an organization.

6. Conclusion

This paper addressed the timely challenge of navigating and charting the model atlas of public repos. Taking Hugging Face as a case study, we showed that its atlas has more complex structure than previously thought. In particular, It is a DAGs instead of a tree and is very deep. We presented several use cases of the model atlas, including analyzing model training trends and completing missing model documentation. As the documented atlas is very incomplete, we proposed a method for using model weights to chart it. It leverages temporal dynamics and real-world priors such as fans and snakes. Extensive in-the-wild experiments demonstrate that our method outperforms previous approaches.

References

- [1] Réka Albert and Albert-László Barabási. Statistical mechanics of complex networks. *Reviews of modern physics*, 74(1): 47, 2002. 7
- [2] Maor Ashkenazi, Zohar Rimon, Ron Vainshtein, Shir Levi, Elad Richardson, Pinchas Mintz, and Eran Treister. Nern-learning neural representations for neural networks. *arXiv preprint arXiv:2212.13554*, 2022. 2
- [3] Omri Avrahami, Dani Lischinski, and Ohad Fried. Gan cocktail: mixing gans without dataset access. In *European Conference on Computer Vision*, pages 205–221. Springer, 2022. 2
- [4] Mathieu Bastian, Sebastien Heymann, and Mathieu Jacomy. Gephi: an open source software for exploring and manipulating networks. In *Third international AAAI conference on weblogs and social media*, 2009. 11
- [5] Nicholas Carlini, Daniel Paleka, Krishnamurthy Dj Dvijotham, Thomas Steinke, Jonathan Hayase, A Feder Cooper, Katherine Lee, Matthew Jagielski, Milad Nasr, Arthur Conmy, et al. Stealing part of a production language model. *arXiv preprint arXiv:2403.06634*, 2024. 2
- [6] Jialuo Chen, Jingyi Wang, Tinglan Peng, Youcheng Sun, Peng Cheng, Shouling Ji, Xingjun Ma, Bo Li, and Dawn Song. Copy, right? a testing framework for copyright protection of deep learning models. In *2022 IEEE symposium on security and privacy (SP)*, pages 824–841. IEEE, 2022. 2
- [7] Derek J. de Solla Price. Networks of scientific papers. *Science*, 149(3683):510–515, 1965. 7, 8
- [8] Amil Dravid, Yossi Gandelsman, Kuan-Chieh Wang, Rameen Abdal, Gordon Wetzstein, Alexei A Efros, and Kfir Aberman. Interpreting the weight space of customized diffusion models. *arXiv preprint arXiv:2406.09413*, 2024. 2
- [9] Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, et al. The llama 3 herd of models. *arXiv preprint arXiv:2407.21783*, 2024. 2, 4
- [10] Gabriel Eilertsen, Daniel Jönsson, Timo Ropinski, Jonas Unger, and Anders Ynnerman. Classifying the classifier: dissecting the weight space of neural networks. *arXiv:2002.05688*, 2020. 2
- [11] Ziya Erkoç, Fangchang Ma, Qi Shan, Matthias Nießner, and Angela Dai. Hyperdiffusion: Generating implicit neural fields with weight-space diffusion. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 14300–14310, 2023. 2
- [12] Sarah Gao and Andrew Kean Gao. On the origin of llms: An evolutionary tree and graph for 15,821 large language models. *arXiv preprint arXiv:2307.09793*, 2023. 2
- [13] Almog Gueta, Elad Venezian, Colin Raffel, Noam Slonim, Yoav Katz, and Leshem Choshen. Knowledge is a region in weight space for fine-tuned language models. In *The 2023 Conference on Empirical Methods in Natural Language Processing*, 2023. 2, 5, 8
- [14] David Ha, Andrew Dai, and Quoc V Le. Hypernetworks. *arXiv preprint arXiv:1609.09106*, 2016. 2
- [15] Vincent Herrmann, Francesco Faccio, and Jürgen Schmidhuber. Learning useful representations of recurrent neural network weight matrices. In *Forty-first International Conference on Machine Learning*, 2024. 2
- [16] Eliahu Horwitz, Bar Cavia, Jonathan Kahana, and Yedid Hoshen. Representing model weights with language using tree experts. *arXiv preprint arXiv:2410.13569*, 2024. 2, 5, 8
- [17] Eliahu Horwitz, Jonathan Kahana, and Yedid Hoshen. Recovering the pre-fine-tuning weights of generative models. In *Forty-first International Conference on Machine Learning*, 2024. 2, 5
- [18] Eliahu Horwitz, Asaf Shul, and Yedid Hoshen. Unsupervised model tree heritage recovery. In *The Thirteenth International Conference on Learning Representations*, 2025. 2, 5, 7, 8
- [19] Gabriel Ilharco, Marco Tulio Ribeiro, Mitchell Wortsman, Suchin Gururangan, Ludwig Schmidt, Hannaneh Hajishirzi, and Ali Farhadi. Editing models with task arithmetic. *arXiv preprint arXiv:2212.04089*, 2022. 8
- [20] Mathieu Jacomy, Tommaso Venturini, Sebastien Heymann, and Mathieu Bastian. Forceatlas2, a continuous graph layout algorithm for handy network visualization designed for the gephi software. *PloS one*, 9(6):e98679, 2014. 1
- [21] Jonathan Kahana, Eliahu Horwitz, Imri Shuval, and Yedid Hoshen. Deep linear probe generators for weight space learning. In *The Thirteenth International Conference on Learning Representations*, 2025. 2
- [22] Jonathan Kahana, Or Nathan, Eliahu Horwitz, and Yedid Hoshen. Can this model also recognize dogs? zero-shot model search from weights. *arXiv preprint arXiv:2502.09619*, 2025. 2, 5
- [23] Miltiadis Kofinas, Boris Knyazev, Yan Zhang, Yunlu Chen, Gertjan J Burghouts, Efstratios Gavves, Cees GM Snoek, and David W Zhang. Graph neural networks for learning equivariant representations of neural networks. *arXiv preprint arXiv:2403.12143*, 2024. 2
- [24] Black Forest Labs. Flux. <https://github.com/black-forest-labs/flux>, 2024. 3
- [25] Derek Lim, Haggai Maron, Marc T Law, Jonathan Lorraine, and James Lucas. Graph metanetworks for processing diverse neural architectures. *arXiv preprint arXiv:2312.04501*, 2023. 2
- [26] Daohan Lu, Sheng-Yu Wang, Nupur Kumari, Rohan Agarwal, Mia Tang, David Bau, and Jun-Yan Zhu. Content-based search for deep generative models. In *SIGGRAPH Asia 2023 Conference Papers*, pages 1–12, 2023. 5
- [27] Michael Luo, Justin Wong, Brandon Trabucco, Yanping Huang, Joseph E Gonzalez, Ruslan Salakhutdinov, Ion Stoica, et al. Stylus: Automatic adapter selection for diffusion models. *Advances in Neural Information Processing Systems*, 37:32888–32915, 2024. 5
- [28] Margaret Mitchell, Simone Wu, Andrew Zaldívar, Parker Barnes, Lucy Vasserman, Ben Hutchinson, Elena Spitzer, Inioluwa Deborah Raji, and Timnit Gebru. Model cards for model reporting. In *Proceedings of the conference on fairness, accountability, and transparency*, pages 220–229, 2019. 2
- [29] Aviv Navon, Aviv Shamsian, Idan Achituve, Ethan Fetaya, Gal Chechik, and Haggai Maron. Equivariant architectures

- for learning in deep weight spaces. In *International Conference on Machine Learning*, pages 25790–25816. PMLR, 2023. 2
- [30] Aviv Navon, Aviv Shamsian, Ethan Fetaya, Gal Chechik, Nadav Dym, and Haggai Maron. Equivariant deep weight space alignment. *arXiv preprint arXiv:2310.13397*, 2023. 2
- [31] Caelean Osborne, Jennifer Ding, and Hannah Rose Kirk. The ai community building the future? a quantitative analysis of development activity on hugging face hub. *Journal of Computational Social Science*, 7(2):2067–2105, 2024. 2
- [32] William Peebles, Ilija Radosavovic, Tim Brooks, Alexei A Efros, and Jitendra Malik. Learning to learn with generative models of neural network checkpoints. *arXiv preprint arXiv:2209.12892*, 2022. 2
- [33] Theo Putterman, Derek Lim, Yoav Gelberg, Stefanie Jegelka, and Haggai Maron. Learning on lorax: G-equivariant processing of low-rank weight spaces for large finetuned models. *arXiv preprint arXiv:2410.04207*, 2024. 2
- [34] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-resolution image synthesis with latent diffusion models. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 10684–10695, 2022. 2
- [35] Mohammad Salama, Jonathan Kahana, Eliahu Horwitz, and Yedid Hoshen. Dataset size recovery from lora weights. *arXiv preprint arXiv:2406.19395*, 2024. 2
- [36] Konstantin Schürholt, Dimche Kostadinov, and Damian Borth. Self-supervised representation learning on neural network weights for model characteristic prediction. *Advances in Neural Information Processing Systems*, 34:16481–16493, 2021. 2
- [37] Konstantin Schürholt, Boris Knyazev, Xavier Giró-i Nieto, and Damian Borth. Hyper-representations as generative models: Sampling unseen neural network weights. *Advances in Neural Information Processing Systems*, 35:27906–27920, 2022.
- [38] Konstantin Schürholt, Giorgos Bouritsas, Eliahu Horwitz, Derek Lim, Yoav Gelberg, Bo Zhao, Allan Zhou, Damian Borth, and Stefanie Jegelka. Neural network weights as a new data modality. In *ICLR 2025 Workshop Proposals*, 2024. 2
- [39] Konstantin Schürholt, Michael W. Mahoney, and Damian Borth. Towards scalable and versatile weight space learning. In *Forty-first International Conference on Machine Learning*, 2024. 2
- [40] Viraj Shah, Nataniel Ruiz, Forrester Cole, Erika Lu, Svetlana Lazebnik, Yuanzhen Li, and Varun Jampani. Ziplora: Any subject in any style by effectively merging lorax. *arXiv preprint arXiv:2311.13600*, 2023. 2, 4
- [41] George Stoica, Pratik Ramesh, Boglarka Ecsedi, Leshem Choshen, and Judy Hoffman. Model merging with svd to tie the knots. *arXiv preprint arXiv:2410.19735*, 2024. 2, 4
- [42] Thomas Unterthiner, Daniel Keysers, Sylvain Gelly, Olivier Bousquet, and Ilya Tolstikhin. Predicting neural network accuracy from weights. *arXiv preprint arXiv:2002.11448*, 2020. 2
- [43] Mitchell Wortsman, Gabriel Ilharco, Samir Ya Gadre, Rebecca Roelofs, Raphael Gontijo-Lopes, Ari S Morcos, Hongseok Namkoong, Ali Farhadi, Yair Carmon, Simon Kornblith, et al. Model soups: averaging weights of multiple fine-tuned models improves accuracy without increasing inference time. In *International conference on machine learning*, pages 23965–23998. PMLR, 2022. 4
- [44] Xet Team. Lfs analysis - hugging face space. <https://huggingface.co/spaces/xet-team/lfs-analysis>, 2025. Accessed: 2025-03-07. 2
- [45] Jiashu Xu, Fei Wang, Mingyu Derek Ma, Pang Wei Koh, Chaowei Xiao, and Muhao Chen. Instructional fingerprinting of large language models. *arXiv preprint arXiv:2401.12255*, 2024. 2
- [46] Prateek Yadav, Colin Raffel, Mohammed Muqeeth, Lucas Caccia, Haokun Liu, Tianlong Chen, Mohit Bansal, Leshem Choshen, and Alessandro Sordani. A survey on model mo-erging: Recycling and routing among specialized experts for collaborative learning. *arXiv preprint arXiv:2408.07057*, 2024. 4
- [47] Prateek Yadav, Derek Tam, Leshem Choshen, Colin A Raffel, and Mohit Bansal. Ties-merging: Resolving interference when merging models. *Advances in Neural Information Processing Systems*, 36, 2024. 2, 4
- [48] Xinyu Yang, Weixin Liang, and James Zou. Navigating dataset documentations in ai: A large-scale analysis of dataset cards on hugging face. *arXiv preprint arXiv:2401.13822*, 2024. 2
- [49] Zhiguang Yang and Hanzhou Wu. A fingerprint for large language models. *arXiv preprint arXiv:2407.01235*, 2024. 2
- [50] Nicolas Yax, Pierre-Yves Oudeyer, and Stefano Palminteri. PhyloLM: Inferring the phylogeny of large language models and predicting their performances in benchmarks. In *The Thirteenth International Conference on Learning Representations*, 2025. 2, 5
- [51] Runpeng Yu and Xinchao Wang. Neural lineage. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 4797–4807, 2024. 2, 5
- [52] Runpeng Yu and Xinchao Wang. Neural phylogeny: Fine-tuning relationship detection among neural networks. In *The Thirteenth International Conference on Learning Representations*, 2025. 2, 5
- [53] Allan Zhou, Kaien Yang, Yiding Jiang, Kaylee Burns, Winnie Xu, Samuel Sokota, J Zico Kolter, and Chelsea Finn. Neural functional transformers. *Advances in neural information processing systems*, 36, 2024. 2

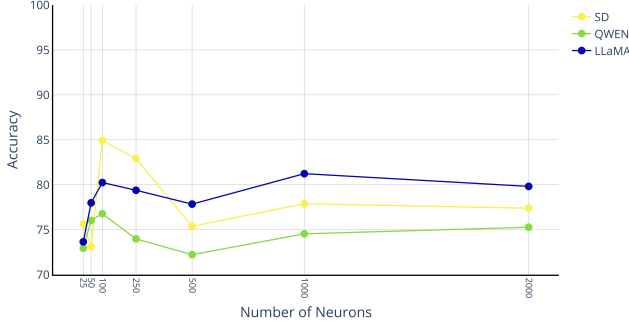


Figure 7. **Number of neurons:** Accuracy as a function of the number of neurons, indeed, 100 presents a good tradeoff between performance and resources.

A. Ablations

In Fig. 7 we ablate how well subsampling the number of neurons affects the model accuracy. We can see that 100 neurons present a good trade off between accuracy and computational cost. In Tab. 4 we provide the runtime of our algorithm for different number of nodes.

B. Additional figures of graphs

In Fig. 13 we visualize the license and `pipeline_tag` of models. In Figs. 10 to 12, 16 and 17 we visualize individual CCs. Figures are compressed to reduce file size.

C. Implementation details

To visualize the atlas, we use the open-source software Gephi [4]. We set $K = 5$ nearest neighbors in all DAGs and determine the snake correlation threshold (ρ_{th}) dynamically for each DAG as the 60th percentile of all node correlations. The distance spreading correlation threshold (K_{th}) is fixed at 0.05.

Additionally, in all our experiments, we use 100 random neurons as model features. For Stable Diffusion, we select these neurons exclusively from the attention layers, as these layers typically undergo LoRA fine-tuning. In all DAGs, we allocate 10% of the actual DAG for training and reserve the remaining 90% as a test set. These hyperparameters are applied consistently across all three tested DAG structure recovery settings.

D. Dataset details

We construct two datasets from the Hugging Face “hub-stats” dataset²: (i) Metadata Imputation and (ii) DAG Recovery. We will make our dataset publicly available on Hugging Face. To process the data, we first used the

²<https://huggingface.co/datasets/cfahlgren1/hub-stats>

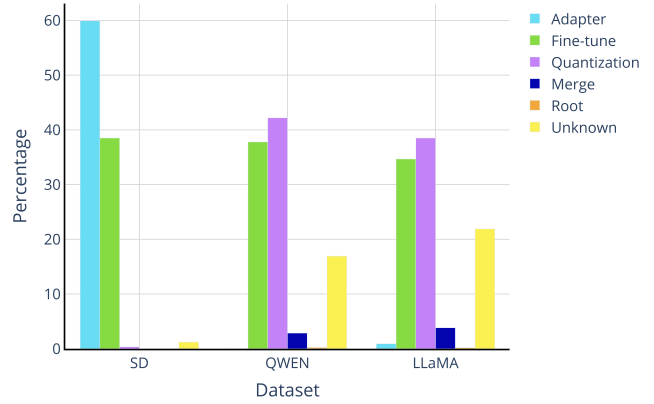


Figure 8. **Relation type in our dataset**

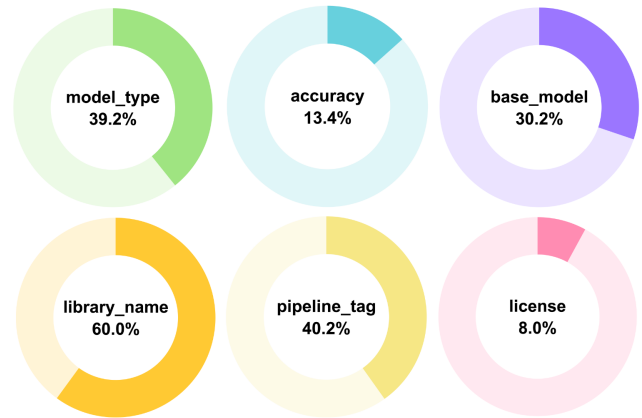


Figure 9. **Documentation level in Hugging Face**

`base_model` attribute to group models into connected components. Since many models lacked this attribute, we manually inspected the largest 800 connected components and filled in the missing values. This step helped merge large components with incomplete information and reduced data noise.

The metadata imputation dataset includes all 1.3 million models. The tasks involve imputing the following attributes: `pipeline_tag`, `library_name`, `model_type`, `license`, and `relation_type`, as well as estimating model evaluation metrics. To obtain the evaluation metrics, we downloaded model cards and extracted the metrics using *deepseek-ai/DeepSeek-R1-Distill-Qwen-7B* in a zero-shot setting on models from *mistralai/Mistral-7B-v0.1* CC.

To construct the charting dataset, we created a graph of all models on Hugging Face. Using the `base_model` attribute, we established ground truth edges between models that reference each other. We removed all nodes with no parent or child (i.e., CCs with a single model), reducing the initial 1.3 million models to approximately 400,000 models

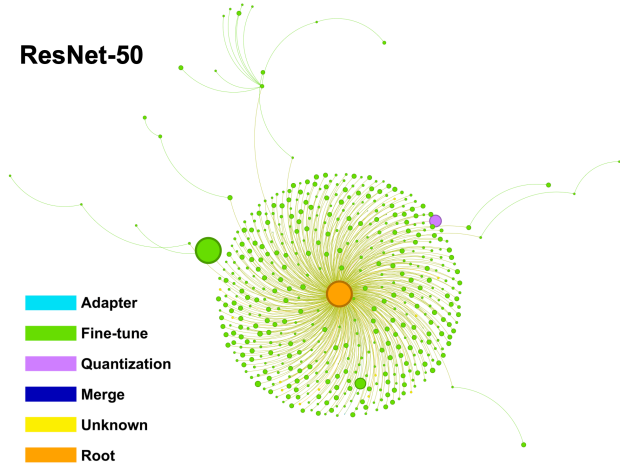


Figure 10. *Individual connected components - 1/5*

Table 4. *Runtime*: Our method can recover an atlas with thousands of models in a matter of seconds.

# of Nodes	Time (sec)
20	0.01
100	0.06
1000	0.58
10000	19.1

across different connected components. We defined sources as nodes with no incoming edges and considered weakly connected components, ensuring that each CC has a single source node.

To facilitate model downloads, we pruned the CCs by randomly sampling K leaf nodes if a node had more than K leaf nodes, filtering out the rest. This allowed us to download hundreds rather than tens of thousands of models. We set $K = 3$. If a download failed, we removed the model from the ground truth atlas.

During the download process, quantized model weights required special handling to convert them back into workable weight files. The conversion was based on the tensor type, which indicates whether a model has undergone quantization. While reviewing the Hugging Face model cards, we discovered multiple tagging inconsistencies—some models were incorrectly labeled as quantized, while others were missing the label. We propose using the tensor type as a reliable indicator of quantization status.

For a breakdown of relation types in our dataset, see Fig. 8.

E. Error analysis and visualizing

In Figs. 14 and 15 we breakdown the errors of our method for all CCs.

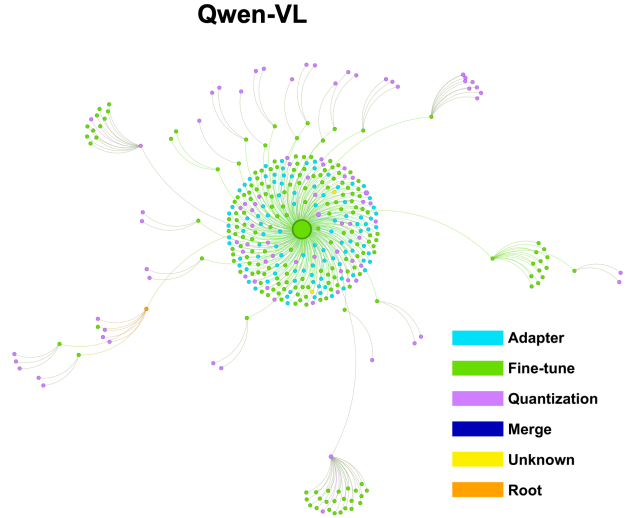


Figure 11. *Individual connected components - 2/5*

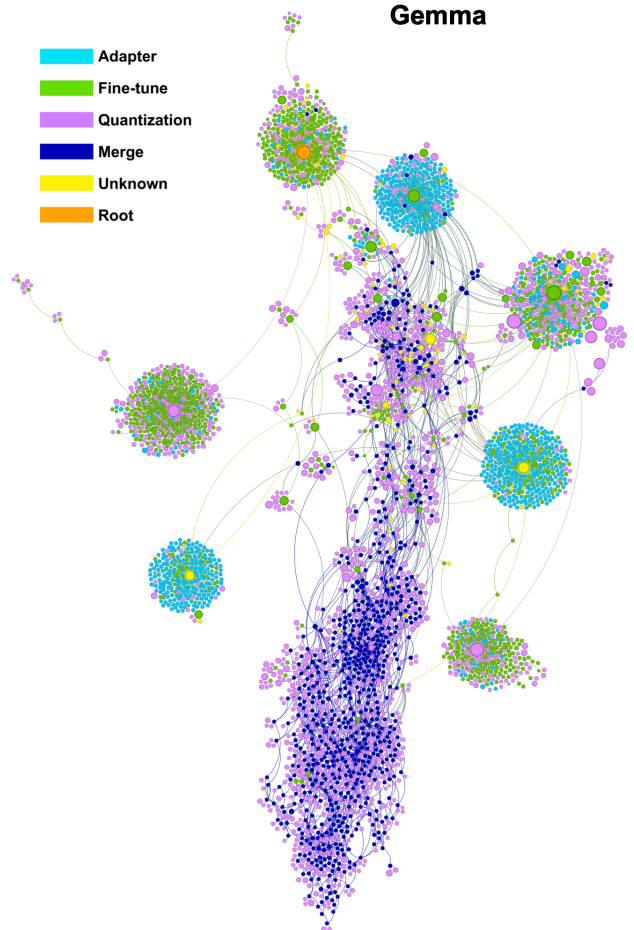
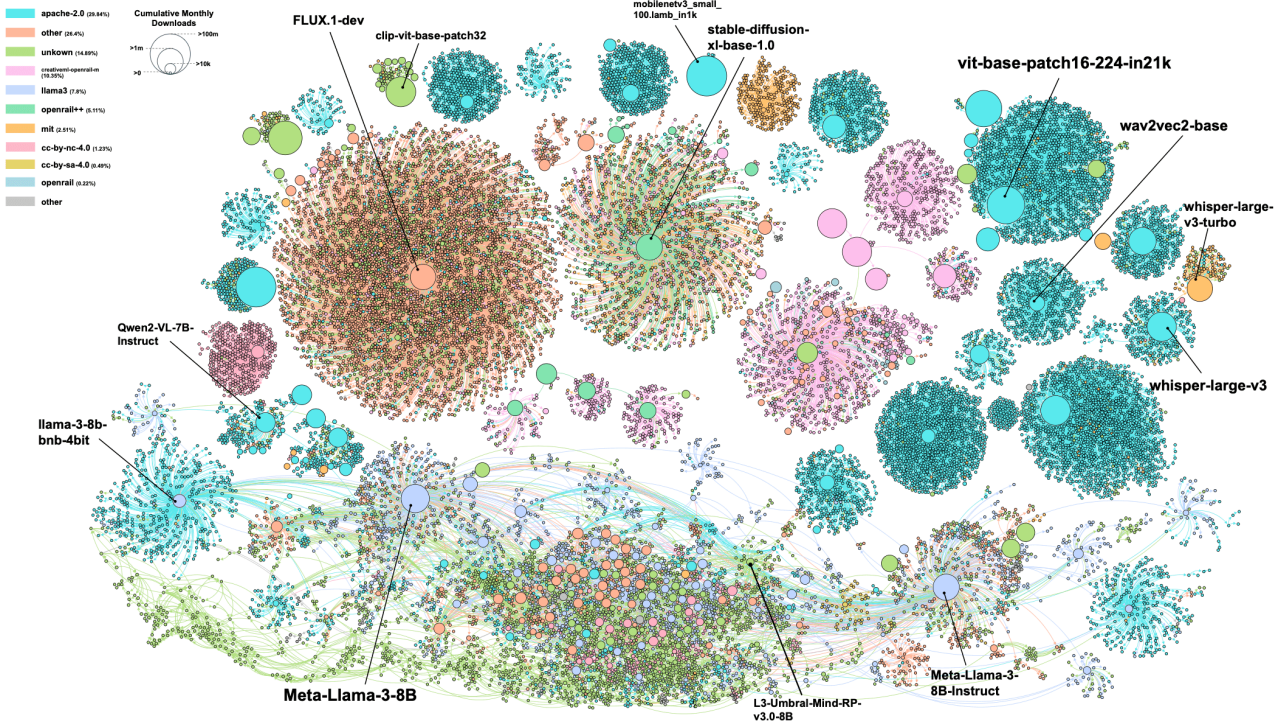


Figure 12. *Individual connected components - 3/5*

License



Pipeline

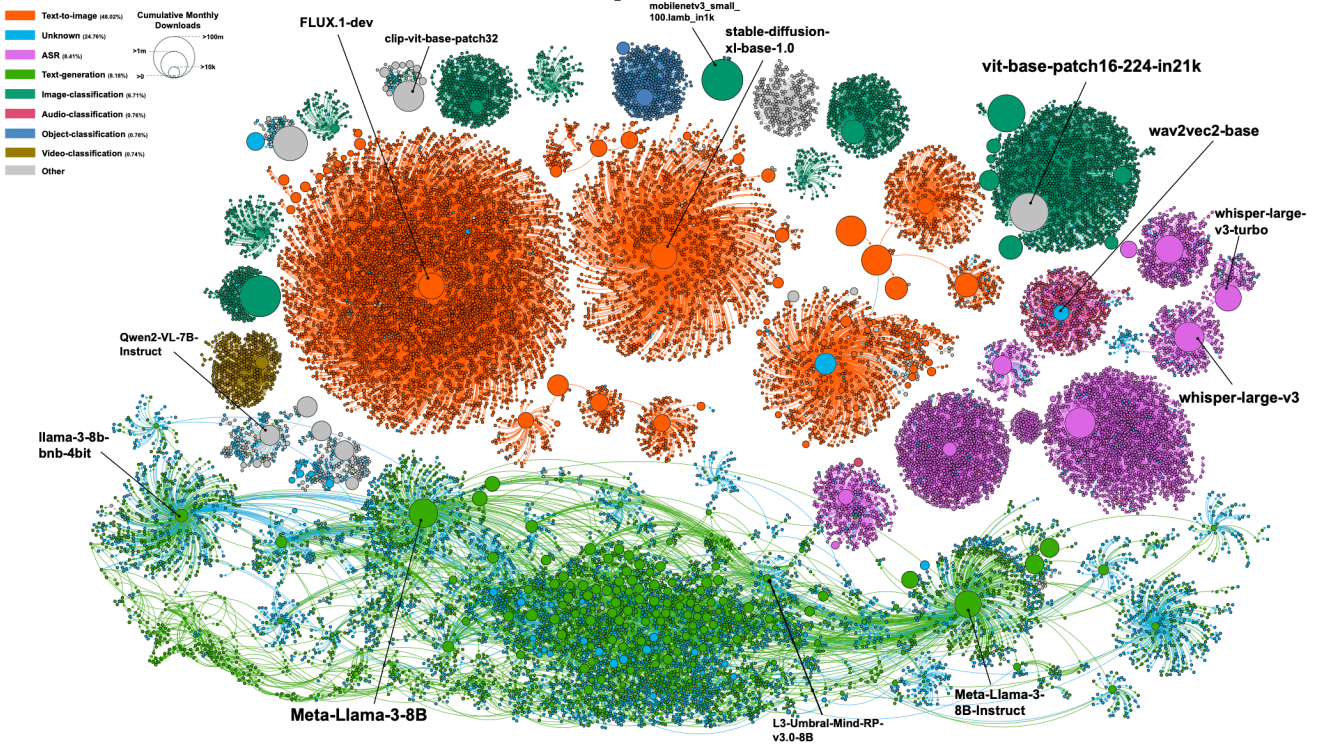


Figure 13. Visualizing other attributes

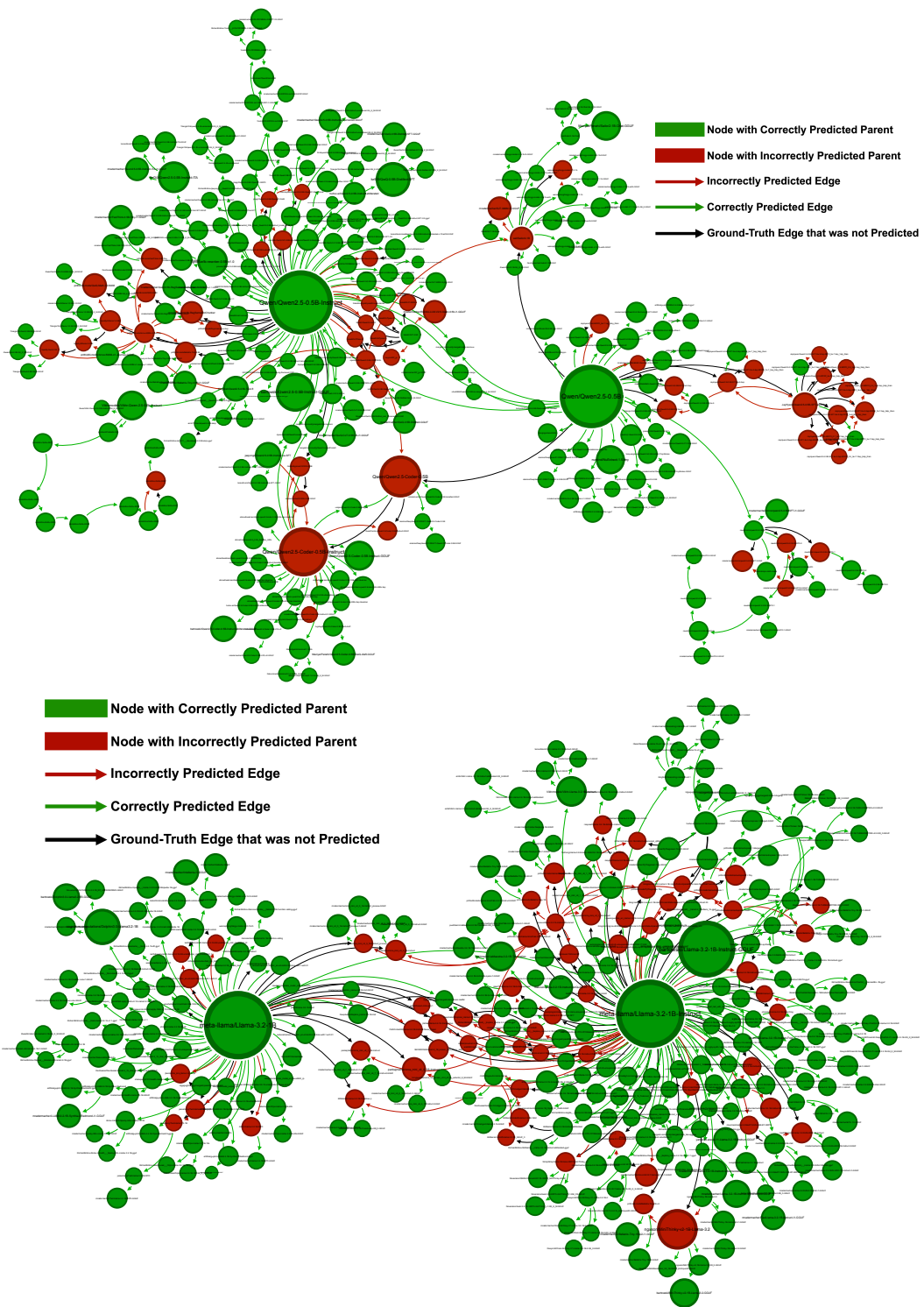


Figure 14. *Atlas prediction errors*: Node names are visible by zooming-in. Correctly predicted edges and their target nodes are marked in green. Incorrect edge predictions and their target nodes are marked in red. Ground-truth edges that were not predicted are marked in black.

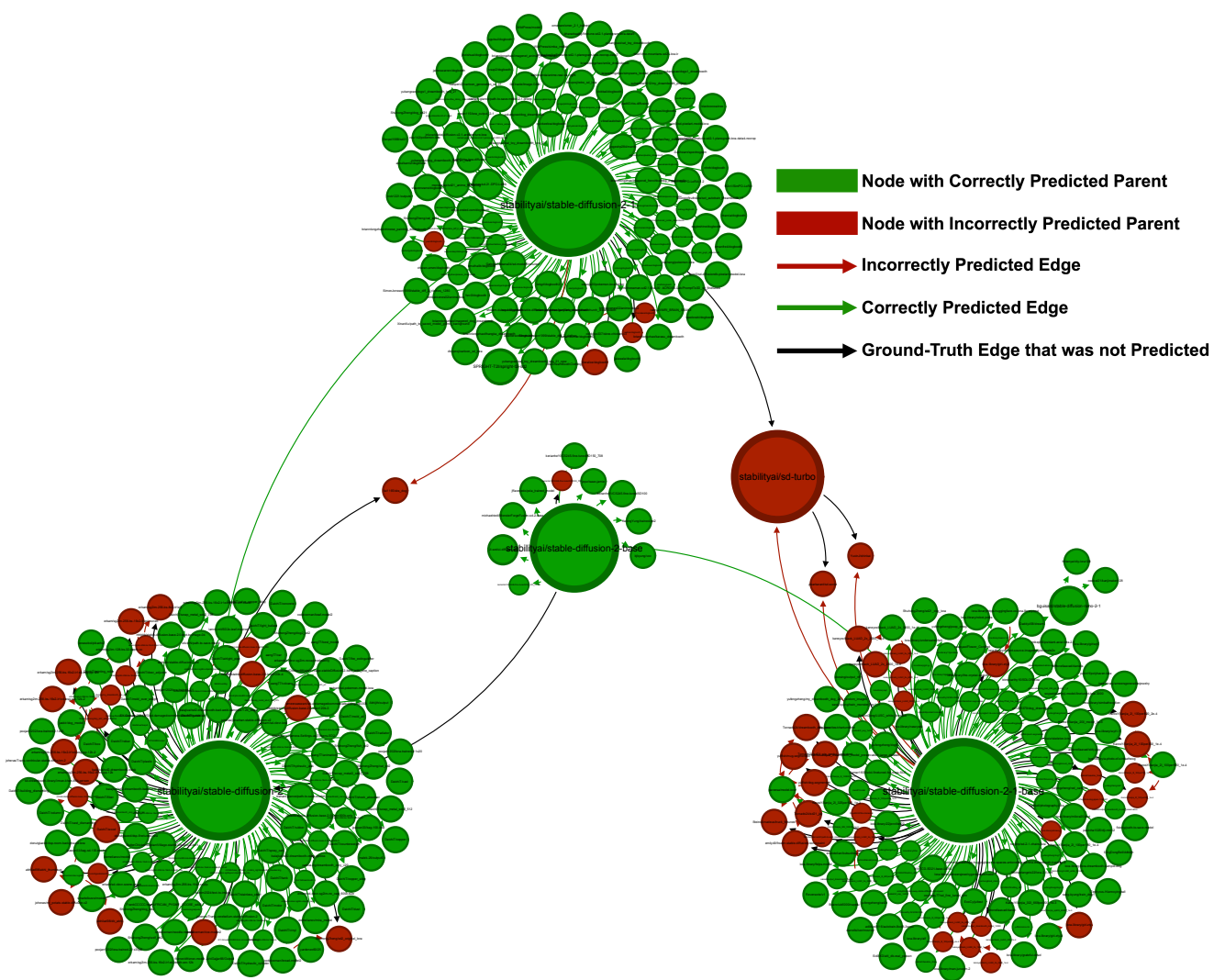


Figure 15. *Atlas prediction errors*: Node names are visible by zooming-in. Correctly predicted edges and their target nodes are marked in green. Incorrect edge predictions and their target nodes are marked in red. Ground-truth edges that were not predicted are marked in black.

Mistral

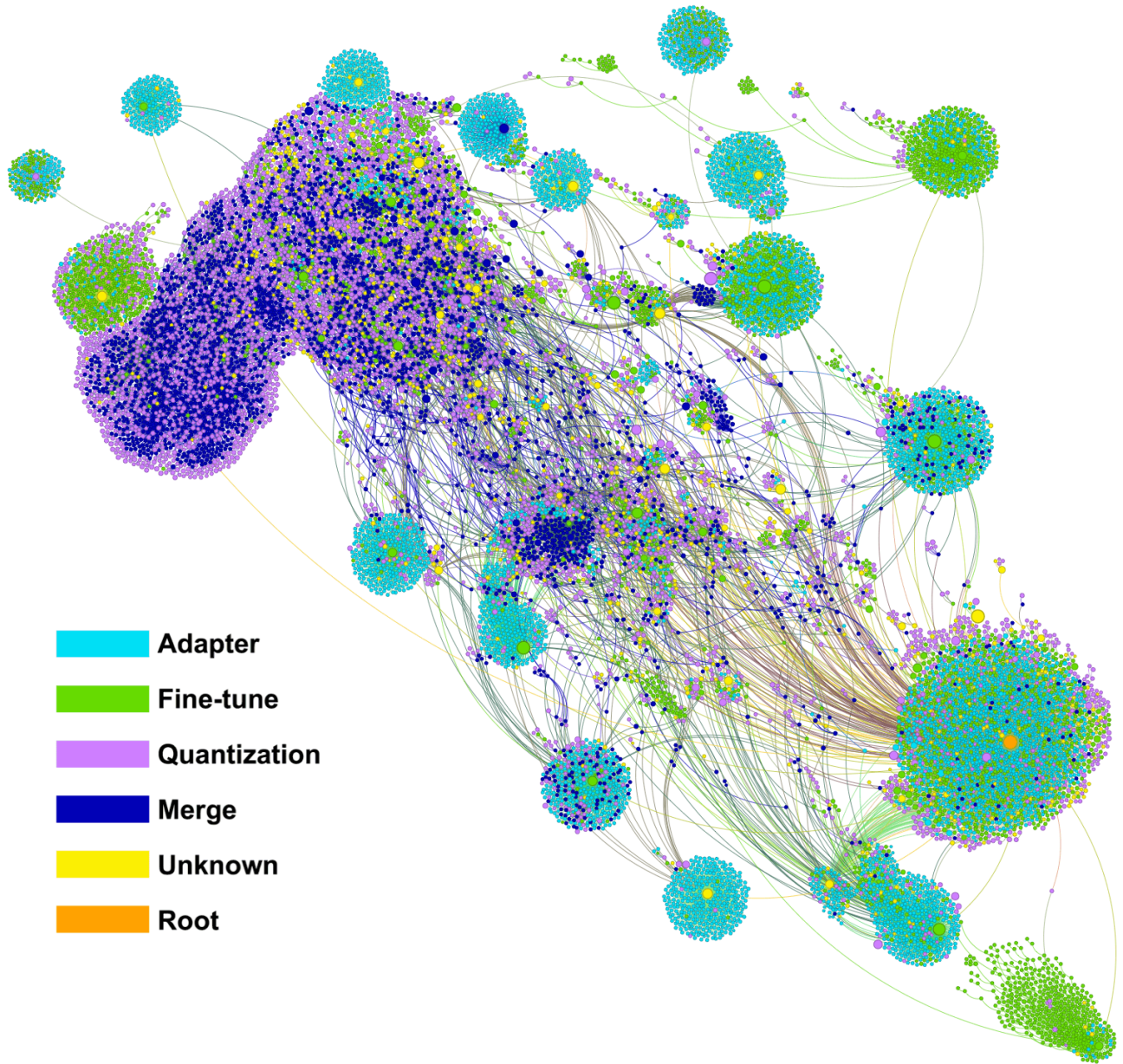


Figure 16. *Individual connected components - 3/5*

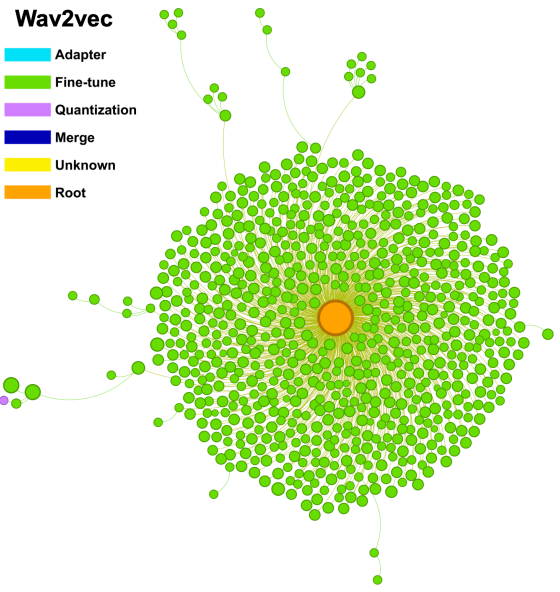
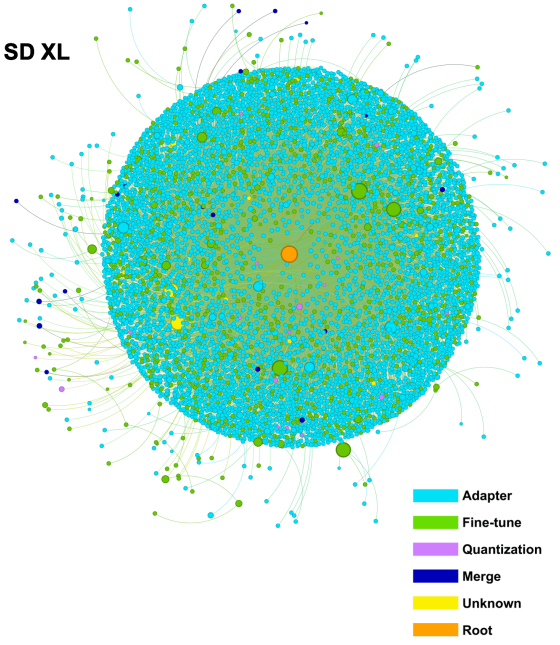
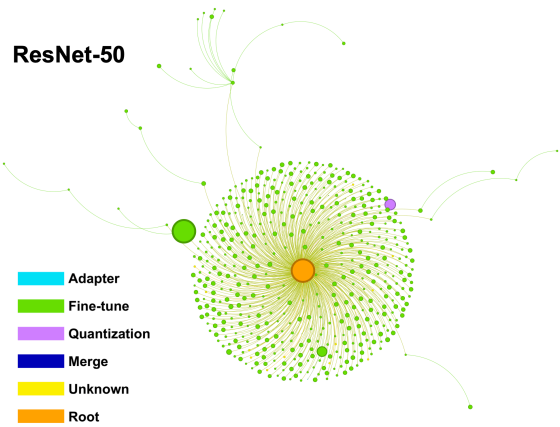


Figure 17. *Individual connected components - 5/5*