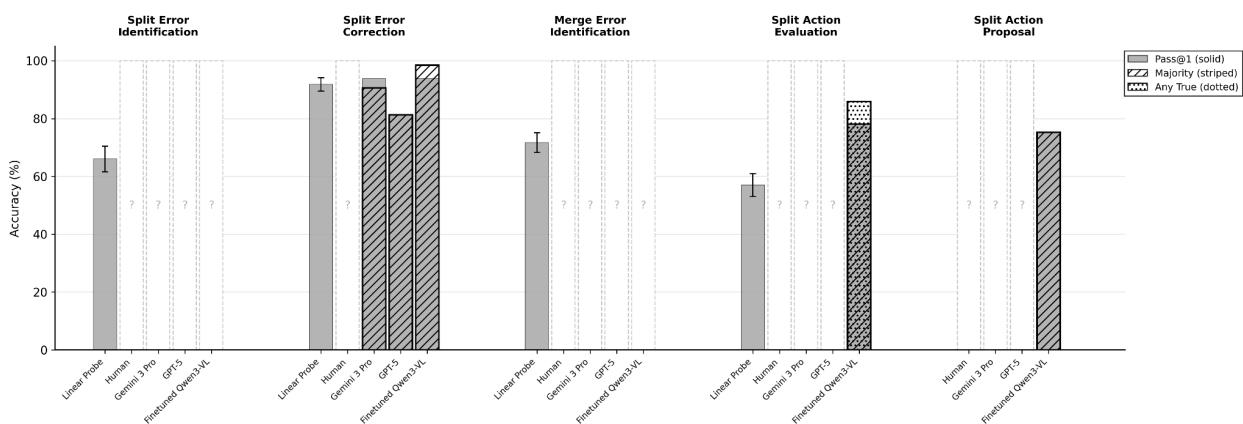


Figure Planning

Human-level proofreading:

- Communicate proofreading subtasks
 - a) Split error identification
 - b) Split error localization
 - c) Split error correction
 - d) Merge error identification
 - e) Split action evaluation
 - f) Split action proposal
 - Examples of each $\frac{1}{3}$ panel (Figure 1)
- Show system for tying it all together
 1. Error candidate generation using the skeleton
 2. Flowchart $\frac{1}{3}$ panel (Figure 1)
- Show human-level performance on the proofreading subtasks (compare against linear probe, human, Gemini 3, gpt-5)
 - Bar charts for a) c) d) e) f)



- CDF for b)
- Figure 2
- Show overall performance
 - Across 100 microns roots (not in ANY of the training set). Table, split by Split errors and Merge errors.

	Split Errors Corrected	Merge Errors Corrected	Split Errors Added	Merge errors added	Errors missed
Human					
Finetuned Model					
NEURD (During)					

rebuttal period)					
ROBOEM (During rebuttal period)					

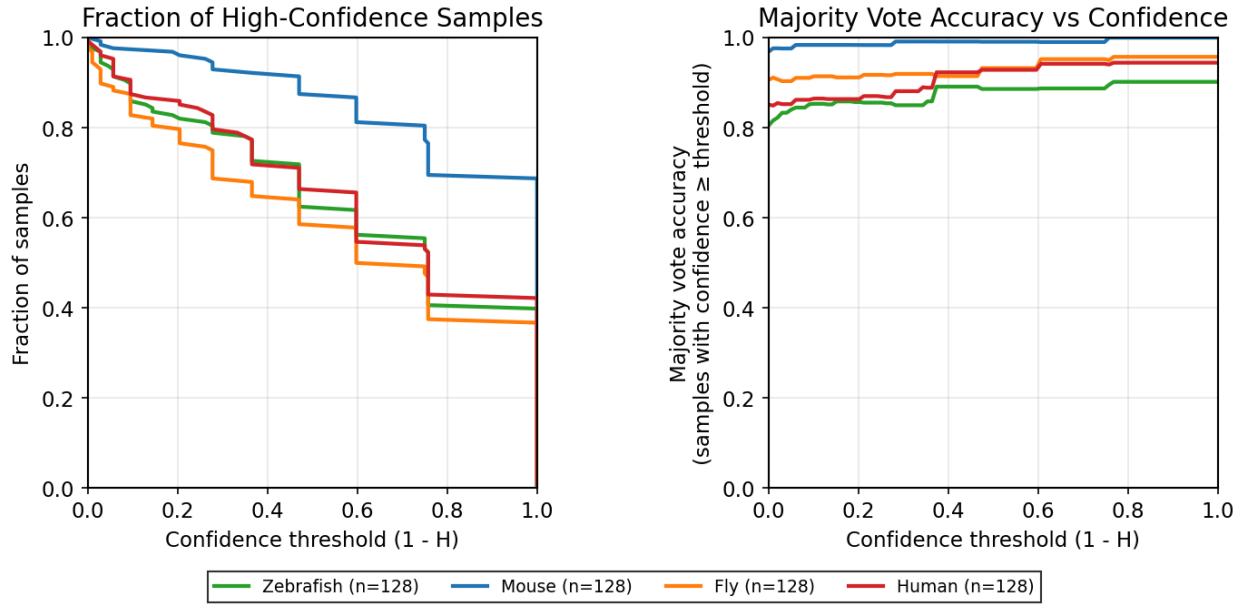
Generalization to new datasets

- Proofreading subtask generalization to Fly, Zebrafish, Mouse
 - Bar chart without any retraining (Figure 3)
- Overall performance (Fly, Zebrafish, Mouse)

	Split Errors Corrected	Merge Errors Corrected	Split Errors Added	Merge errors added	Errors missed
Fly					
Zebrafish					
Human					

Calibrated behavior

- Relationship between sample variance and accuracy

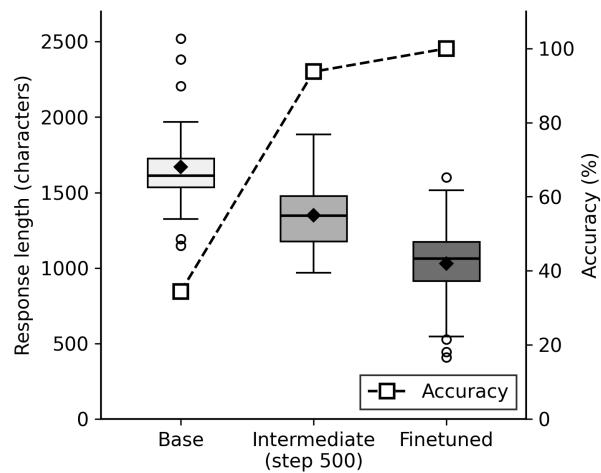


- Can the model faithfully express confidence?
- Does it faithfully describe what it sees?

Computational costs

Interpretability

- What does the model look at in the image vs the linear probe?
- How does the behavior change



- What are the rules/heuristics that it learns?

GPT-5	<p>Heuristics GAINED through training:</p> <ul style="list-style-type: none"> - Multi-view alignment rigor: explicit checks for misregistration/offset; uniform shape/size profile across the boundary; consistency "across color change" - Depth/volumetric reasoning: explicit over/under pass detection; "threading through pores" rejection; insistence on a single volumetric junction point and tapering/blending at true joins - Morphology/texture consistency: matching surface texture (e.g., varicosities) and morphology class (thin filament vs thick/branched) across the junction - Branch-pattern continuity: not just "a branch exists," but lateral branches and overall branching pattern continue smoothly across the junction <p>Heuristics LOST (or de-emphasized):</p> <ul style="list-style-type: none"> - Naive branch cue: early acceptance of generic T/Y-junctions as positive is replaced by stricter branch-pattern continuity - "Parallel run-up" as a positive cue is tempered—parallel adjacency without true convergence is now a reject - Less reliance on single-view/projection cues; simple colinearity without volumetric confirmation is no longer sufficient <p>What explains the accuracy improvement (34% → 100%):</p> <ul style="list-style-type: none"> - Addition of depth- and volume-aware negatives (over/under, threading, parallel adjacency, butt joins without blending) eliminates common false positives from projection artifacts - Stronger multi-cue agreement requirements (alignment, curvature, thickness, texture, branch-pattern continuity) reduce ambiguous cases and enforce true anatomical continuity - Morphology-class and texture checks prevent mismatched merges (thin vs thick/branched), improving both precision and recall
Gemini 3 Pro	<p>What heuristics are GAINED through training?</p> <ul style="list-style-type: none"> - Morphological Semantics: The model learns to look beyond simple geometry to distinct structural textures and complexity (e.g., rejecting merges between "smooth/linear" and "brambly/branching" segments), which appears first in the Intermediate checkpoint - 3D Relational Context: The reasoning shifts from simple "alignment" to understanding complex spatial behaviors like "weaving through," "passing over," or "crossing" (Trajectory Independence), distinguishing true connections from incidental spatial overlaps <p>What heuristics are LOST?</p> <ul style="list-style-type: none"> - Micro-Boundary Detection: The Base model relies on "Intact Boundaries" and specific "membranes" to reject merges. Later models drop this specific edge-detection focus in favor of broader morphological consistency - Projection-based Hesitation: The Base model treats "Proximity Without Contact" as a specific visual heuristic (likely confused by 2D projections). Later models incorporate this directly into 3D spatial reasoning ("Spatial Separation" in all views) rather than treating it as a separate adjacency issue <p>What explains the accuracy improvement (34% → 100%)?</p>

	- From Geometry to Topology: The Base model relies on local geometric cues (Do they touch? Do they line up?), which leads to frequent errors where distinct neurons simply cross paths or touch. The Final model applies semantic consistency; it correctly identifies that if one segment is smooth and the other is complex/branching, they cannot represent the same biological structure, even if they physically align or touch
--	--

Cost

- Retrospective calculation of the cost of proofreading MICRONS, FlyWire
- Forward analysis for cubic millimeter LICONN/barcoding?

Appendix

	Identification			Correction			
	Recall	Precision	FPR	Correct	Wrong	None when Should	None correct
Human							
Finetuned Model							
Gemini 3							
GPT-5							

Tab 1

ConnectomeAgent

What do we want out of AI proofreaders?

- First, we want a human-level of performance on proofreading tasks: split error identification, split error correction, merge error identification, and merge error correction.
- Second, we want calibrated behavior. The models should seek more information when they are unsure about what do.
- Third, we want proofread performance within dataset, and cross-dataset/cross-species generalization without a large amount of retraining.
- Fourth, we want faithful descriptions of what is being seen/we want models to get the answers right for the right reasons. (optional)

Human-level performance on proofreading tasks

- a) We need methods for identifying potential split errors and merge errors

i) **Split errors:**

[✓] Todo: recall, precision of split error @ endpoints across large number of neurons in multiple datasets (fly, mouse, human, fish) glebrazgar@gmail.com

- **Results:** 96,7% split error recall, 2% precision.
- **Methods:**
 - Localization:
 1. Extract the skeleton and find nodes with 1 edge (endpoints).
 - Data generation (labelling):
 2. Get merge corrections from history and extract split error locations.
 3. Match endpoints within 2µm of corrections = positive labels.
 4. Render 3-view images at each endpoint.

ii) **Merge errors (WIP)**

Todo: recall, precision of merge error at junctions across large number of neurons in multiple datasets (mouse, human, fish, fly in that order)

glebrazgar@gmail.com

- **Results:** 85% merge error recall, 2% precision.
- **Methods Old:**
 - Localization:
 1. Get all fragments of the neuron
 2. Get skeleton for each fragment
 3. Add bridge edges using error locations (effectively re-constructing the old merged neuron from its fragments)

4. Find branchpoints on the combined graph

- b) We need performance metrics for error identification. This is WITHOUT correction, just identification
- i) This best-case scenario is the perfect recall with maximum allowed precision, since this is a filter.
 - ii) endpoint error candidates: Is there a split error here, yes or no?
 - 1) TODO: Generate large datasets (start with 512, then scale to O(10k), balanced TP/TN) for mouse, human, fish, fly in that order using the question dataset format. @gleb
 - 2) TODO: Evaluate recall and precision @jeff (Blocked by bii1)
 - (a) +/- finetuning 32B
 - (b) Best models
 - (i) Gemini 3 Pro
 - (ii) GPT-5
 - (c) Human performance, on subset of 512
 - (d) Linear probe on siglip2-so400m-patch16-512
 - iii) Merge error candidates: Is there a merge error here, yes or no?
 - 1) TODO: Generate large datasets (start with 512, then scale to O(10k), balanced TP/TN) for mouse, human, fish, fly in that order using the question dataset format. @gleb ,
 - 2) TODO: Evaluate recall and precision @jeff (Blocked by biii1)
- c) We need performance metrics for error correction.
- i) For merge action,
 - 1) Binary performance @jeff
 - (a) Mouse

```
7. 5-Fold Cross-Validation (linear):
Fold accuracies: ['76.7%', '76.7%', '82.4%', '80.4%', '78.4%']
Fold precisions: ['81.4%', '79.2%', '85.1%', '75.4%', '76.4%']
Fold recalls: ['68.6%', '73.1%', '78.4%', '90.2%', '82.4%']
Fold F1 scores: ['74.5%', '76.0%', '81.6%', '82.1%', '79.2%']

Mean CV accuracy: 78.9% (+/- 4.4%)
Mean CV precision: 79.5% (+/- 7.0%)
Mean CV recall: 78.5% (+/- 14.9%)
Mean CV F1: 78.7% (+/- 6.1%)
```

(i) -Finetuning Qwen3 VL 32B:

```
"base_model": "Qwen/Qwen3-VL-32B-Instruct",
"dataset_source": "merge-parquet",
"split": "val",
"num_samples": 128,
"num_correct": 41,
"accuracy": 0.3203125,
```

(ii) +finetuning 32B SFT, 4> positive responses from GPT-5 or Gemini 3 Pro

(1) Finetuned with answer only

Prompted to just provide the answer

```
=====
Results
=====
Accuracy: 97.66% (125/128)
=====
Results saved to: /results/eval_merge_action_checkpoints_merge_action_finetune_Qwen3-VL-32B-Instruct_merge_action_32B_big_20260112_143952_samples.json
Final Accuracy: 97.66%
```

Prompted to provide the answer and the analysis

```
=====
Results
=====
Accuracy: 92.19% (118/128)
=====
Results saved to: /results/eval_merge_action_checkpoints_merge_action_finetune_Qwen3-VL-32B-Instruct_rationale_dropout_ochs1_lr0.0002_r16_merged_val_20260109_221028.json
```

(2) Answer with traces (consistently performs worse)

(iii) Best models,

(1) Gemini 3

```
=====
Results
=====
Individual Response Accuracy: 87.66% (561/640)
Majority Vote Accuracy: 90.62% (116/128)
=====
Results saved to: evaluation_results/eval_merge_action_google_gemini-3-pro-preview_votes5_20260110_125432.json
```

(2) GPT-5

```
=====
Results
=====
Individual Response Accuracy: 80.16% (513/640)
Majority Vote Accuracy: 81.25% (104/128)
=====
Results saved to: evaluation_results/eval_merge_action_gpt-5_votes5_20260110_124127.json
```

(iv) Human performance, on subset of 512

(v) Linear probe on siglip2-so400m-patch16-512

```
7. 5-Fold Cross-Validation (linear):
Fold accuracies: ['93.2%', '93.2%', '87.3%', '94.1%', '92.2%']
Fold precisions: ['91.1%', '92.6%', '95.5%', '92.7%', '94.0%']
Fold recalls: ['96.2%', '94.3%', '79.2%', '96.2%', '90.4%']
Fold F1 scores: ['93.6%', '93.5%', '86.6%', '94.4%', '92.2%']

Mean CV accuracy: 92.0% (+/- 4.9%)
Mean CV precision: 93.2% (+/- 2.9%)
Mean CV recall: 91.3% (+/- 12.8%)
Mean CV F1: 92.0% (+/- 5.6%)
```

(vi) MLP

```
7. 5-Fold Cross-Validation (mlp):
  Fold accuracies: ['91.3%', '93.2%', '86.3%', '89.2%', '75.5%']
  Fold precisions: ['89.3%', '91.1%', '95.3%', '85.0%', '74.5%']
  Fold recalls: ['94.3%', '96.2%', '77.4%', '96.2%', '78.8%']
  Fold F1 scores: ['91.7%', '93.6%', '85.4%', '90.3%', '76.6%']

  Mean CV accuracy: 87.1% (+/- 12.5%)
  Mean CV precision: 87.1% (+/- 14.2%)
  Mean CV recall: 88.6% (+/- 17.2%)
  Mean CV F1: 87.5% (+/- 12.2%)
```

- 2) TODO: Multiple choice of N nearest segments (filtered, <= 4, + Refusal),
(a) Before finetuning 32B
(b) After finetuning 32B (working on it)
(c) Best models

ii) For split action

- 1) Binary performance on evaluating good/bad splits
(a) Generating data (@tim)
(b) Linear probe on siglip2-so400m-patch16-512 (1/9/25)

```
7. 5-Fold Cross-Validation (linear):
  Fold accuracies: ['60.2%', '62.1%', '52.9%', '57.8%', '52.0%']
  Fold precisions: ['60.0%', '65.2%', '54.0%', '58.2%', '52.7%']
  Fold recalls: ['67.9%', '56.6%', '51.9%', '61.5%', '55.8%']
  Fold F1 scores: ['63.7%', '60.6%', '52.9%', '59.8%', '54.2%']

  Mean CV accuracy: 57.0% (+/- 8.0%)
  Mean CV precision: 58.0% (+/- 8.9%)
  Mean CV recall: 58.8% (+/- 11.0%)
  Mean CV F1: 58.3% (+/- 8.1%)
```

(c) Gemini

```
=====
Results
=====
Accuracy: 76.56% (98/128)
=====

Results saved to: evaluation_results/eval_split_action_google_gemini-3-pro-preview_20260109_1
```

(d) GPT-5 (pass at 1)

```
=====
Results
=====
Accuracy: 77.34% (99/128)
=====

Results saved to: evaluation_results/eval_split_action_gpt-5_20260110_1
```

- 2) TODO: Performance on generating good/bad splits (less important)

Information seeking to improve performance

Often time, the information needed to make a decision is not available in single snap shot of a scene.

- a) Build an environment where the model can interact freely with the data (Done)
 - i) Move, zoom, rotate

Demo of proofreading within species/dataset

Demo of transfer between species/dataset

VLMs produce faithful descriptions of what is being seen (Optional)

Reliable semantic interventions (Optional)

Notes/Scratch Pad

Proofreading Pipeline

First, classification of segment types (ResNet, open-source VLM, proprietary VLM)

- Single soma
- Multiple somas
- Processes
- Non-neuronal structures
- Nucleus

Two parallel channels;

- Single soma and processes
 - 1. Identify merge action sites (using the skeleton's, most likely)
 - The metric shows that this recovers X% of the split errors in the datasets
 - Likely problem— Wayyy too many candidates
 - 2. Agent corrects splits

- Merge action success/failure rate
- Multiple somas
 1. Idea one: Identify split action branch
 2. Use branching in the skeleton

Overflow

Audience: ICML reviewers. Likely do not know anything about connectomics, and likely know very little about biology. The story perspective is around showing the following results

- Relevance of the problem/Set up
 - High performance
 - Generalization
 - Calibration
 - Interpretation
-

What is connectomics?

Connectomics is focused on developing ground-truth maps of how neurons in the brain are connected.

Why does that matter?

High-resolution maps can be used to simulate nervous systems (Haspel), identify the structural signature of disease, and identify models/value functions in brains relevant to building future AI systems.

What's the problem?

Connectomics maps are generated by imaging the brains of organisms at high resolution ($O(\text{nanometer})$) either using electron microscopy (EM citation, MICRONS, FLYWIRE) or expansion microscopy (LICONN, PRISM citation). Segmentation algorithms (FFN, LSD, citations) are then applied to this data to segment it into the 3D structure of each individual neuron. Due to variations in data quality and organisms, these algorithms make mistakes and require a follow-up proofreading step, in which human annotators review and correct the data. FlyWire quote about requiring 30 human years to proofread.

Multimodal AI systems can do this task

Last year, we published a paper showing that frontier models have sufficiently strong visual priors to perform reasonably well across various proofreading subtasks without any finetuning (ConnectomeBench citation).

Previous work showed that frontier models have sufficiently strong visual priors to perform reasonably well across various proofreading sub-tasks without any finetuning

For example, UNISPAC introduced a system that unified 2D segmentation and 3D automated tracing and addressed the challenge of cross-species generalization by using TTT and SFT to adapt the model to new datasets based on human proofreading corrections.

Why use an image rather than a graph-based approach?

The rise of AI capabilities in visual reasoning suggests that these models can use generalizing visual priors to solve this problem. Vision and reasoning are mechanisms by which people solve this problem. Great if we could get models to do the same things. While mesh- or graph-based approaches are complementary, they do not benefit from the scaling effects of the most potent vision-language models.

Breaking down the major problems:

There are two major classes of errors: split errors and merge errors. Split errors occur when a segment corresponding to one neuron is broken into multiple segments. Merge errors occur when segments from multiple neurons are merged into a single segment. We break proofreading into two phases: identifying errors and correcting them. There is a third component: error localization (i.e., finding potential error candidates in the 3D volume). For this, we leverage heuristics that leverage the segment's skeleton to identify potential error candidates. We find that these heuristics can identify split and merge error locations with ~95% and ~60% recall, respectively, and ~2% precision. *Further methods in the appendix after talking to G.*

Data generation

For each task, we present the models with three views of the visual scene, each from a different orthographic projection, to provide 3D context. For the split error identification problem, we also incorporate 3D slices of the EM or ExM data for additional context. To prevent contamination between training splits, we ensure that the same segments and locations do not appear in different splits. *Full breakdown of the data generation process in the appendix.*

tfarkas@mit.edu Fill in the split generation and evaluation

Training Methods

Model Training and Evaluation

Vision-Language Model Fine-tuning

We fine-tune Qwen3-VL-32B-Instruct, a 32-billion parameter vision-language model, using Low-Rank Adaptation (LoRA) with rank 16. Training is performed on Modal cloud infrastructure using 2x H100 GPUs with the Unsloth library for memory-efficient optimization. LoRA adapters are applied to the language model's attention layers (query, key, value, and output projections) and MLP layers (gate, up, and down projections), as well as the vision-language merger module that bridges the vision encoder to the language model. The vision encoder itself remains frozen. We use the AdamW optimizer with 8-bit precision, a learning rate of 2e-4 with reduce-on-plateau scheduling, and a maximum of 500 training steps per task.

We train separate LoRA adapters for five proofreading tasks: merge error identification, merge action verification, split action verification, endpoint error identification, and endpoint error identification with EM context. Each task uses a batch size of 4-8 with class balancing to address label imbalance—split action uses undersampling while others use oversampling. All tasks are framed as binary classification problems where the model must output "yes" or "no" within XML answer tags. Training data consists of 3-view orthogonal renderings (front, side, top) of neuronal segments, with stratified train/validation/test splits (128 validation, 512 test samples) using group-based splitting by spatial location to prevent data leakage between splits.

The training pipeline employs lazy image loading to handle large datasets efficiently, loading images on-the-fly during batch collation rather than upfront. The model is trained only on assistant responses (not user prompts) using Unsloth's vision data collator with Qwen3-VL chat template delimiters. Checkpoints are saved every 100 steps with the best model selected by validation loss. Training configurations, test set indices, and dataset hashes are persisted alongside model weights to ensure reproducibility and enable consistent evaluation across runs.

Vision-Language Model Evaluation

Fine-tuned LoRA adapters are evaluated on held-out test sets using FastVisionModel inference with the base Qwen3-VL-32B-Instruct model. For each test sample, the model receives the same 3-view orthogonal renderings used during training along with the task-specific prompt. The processor applies the Qwen3-VL chat template to format inputs, and the model generates responses with a maximum of 512 new tokens using greedy decoding. Answers are extracted from the generated text by parsing XML answer tags (`<answer>yes</answer>` or `<answer>no</answer>`) and compared against ground truth labels using exact string matching.

Evaluation supports several ablation modes: (1) base model evaluation without adapters to measure zero-shot performance, (2) blank image controls where all images are replaced with uniform gray 1024×1024 canvases to test for prompt exploitation, (3) simple prompt controls using generic "What do you see?" prompts to verify task-specific reasoning, and (4) answer-only

mode that removes chain-of-thought analysis instructions to isolate decision-making from reasoning. Test set indices are loaded from `test_indices.json` files saved during training to ensure evaluation uses the exact same held-out samples across runs. Results are saved to parquet files containing predictions, ground truth labels, full prompts, model responses, and sample identifiers for downstream analysis.

CNN Baseline Training

As a non-transformer baseline, we fine-tune pretrained ResNet-50 models (initialized with ImageNet weights) on the same proofreading tasks. For tasks with multiple images per sample, images are arranged into grids using automatic layout detection (e.g., 1×2 for 2 images, 2×2 for 4 images) before being resized to 224×224 pixels—the standard ResNet input size. We apply ImageNet normalization ($\text{mean}=[0.485, 0.456, 0.406]$, $\text{std}=[0.229, 0.224, 0.225]$) and data augmentation during training including random horizontal flips, $\pm 10^\circ$ rotation, and color jittering (brightness and contrast $\pm 20\%$).

Training uses the AdamW optimizer with a learning rate of $1e-3$, weight decay of $1e-4$, and cosine annealing LR scheduling over 10 epochs. We use large batch sizes (256) to leverage the smaller model size compared to VLMs, with 8 dataloader workers for efficient I/O. The final fully-connected layer is replaced to match the number of task classes, and we support two training modes: (1) full fine-tuning where all parameters are trainable, and (2) linear probe mode where the backbone is frozen and only the classifier head is trained. Early stopping with patience of 5 epochs is applied based on validation loss.

Dataset splits use the same stratified group-based splitting as VLM training, with identical train/val/test indices saved to `test_indices.json` for cross-model comparison. Class balancing is applied via oversampling minority classes to match the majority class count. Models are trained on a single A10G GPU using standard PyTorch training loops with cross-entropy loss. The best model checkpoint (by validation accuracy) is saved and used for final test set evaluation.

CNN Baseline Evaluation

ResNet models are evaluated on held-out test sets by loading the best checkpoint and computing per-class accuracy. For each test sample, images are converted to grid format with the same layout and preprocessing used during training. The model outputs logits for each class, and predictions are obtained via argmax. Final metrics include overall accuracy and per-class accuracy to identify class-specific performance differences. Evaluation supports the same test set filtering as VLM evaluation, loading indices from `test_indices.json` to ensure fair comparison. Cross-dataset evaluation is also supported by specifying alternative dataset paths, enabling assessment of model generalization to different species or imaging conditions.

Subtask Breakdown

There are five subtasks:

- Split error identification
 - Views of the endpoint of a neuron mesh. Centered on the skeleton endpoint. With EM data right now. Split error correction
 - Views of a neuron mesh in blue and a potential merge candidate in orange
- Split error correction
- Merge error identification
 - Views of neuron mesh are locations where there are and are not merge errors
- Split correction evaluation
 - Views of how the segment would look after a proposed split was executed
- Split correction proposal
 - Views of how the segment looks before a split with a graph overlay

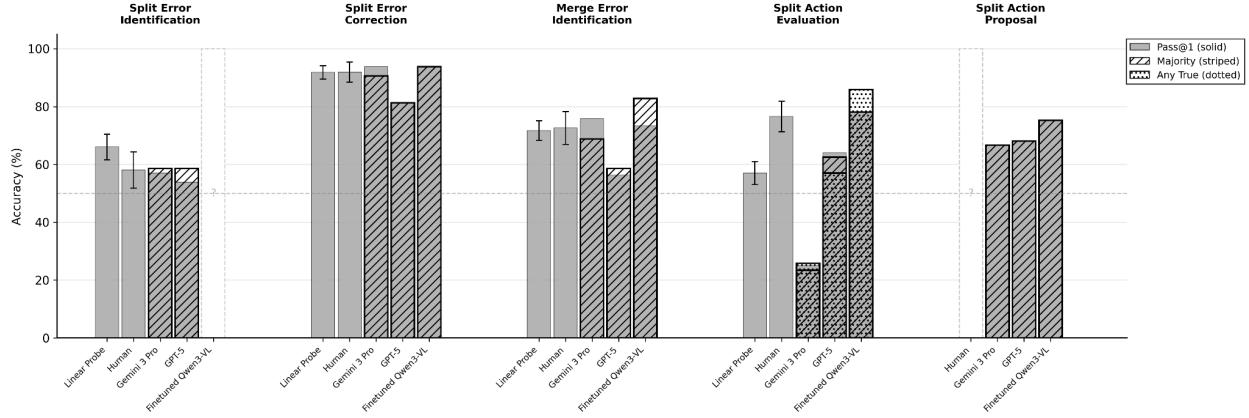
Specify the data distribution in the training splits.

task_name	raw_samples	after_filtering	filtering_applied
segment_classification	HuggingFace	HuggingFace	Exclude classes e/f/g
split_action	36,078	9,542	Dedup by location + undersample
merge_action	13,020	13,020	None
merge_action_multiple_choice	6,301	6,301	None
endpoint_error_identification	2,048	2,048	None
endpoint_error_identification_with_em	1,000	1,000	None
endpoint_localization	2,691	2,691	None
split_proposal	7,060	7,060	None
segment_identity	3,976	3,976	None
merge_error_identification	13,544	11,998	Dedup by (root_id; interface_point)

Results

Performance within distribution

Subtask Performance on MICrONS



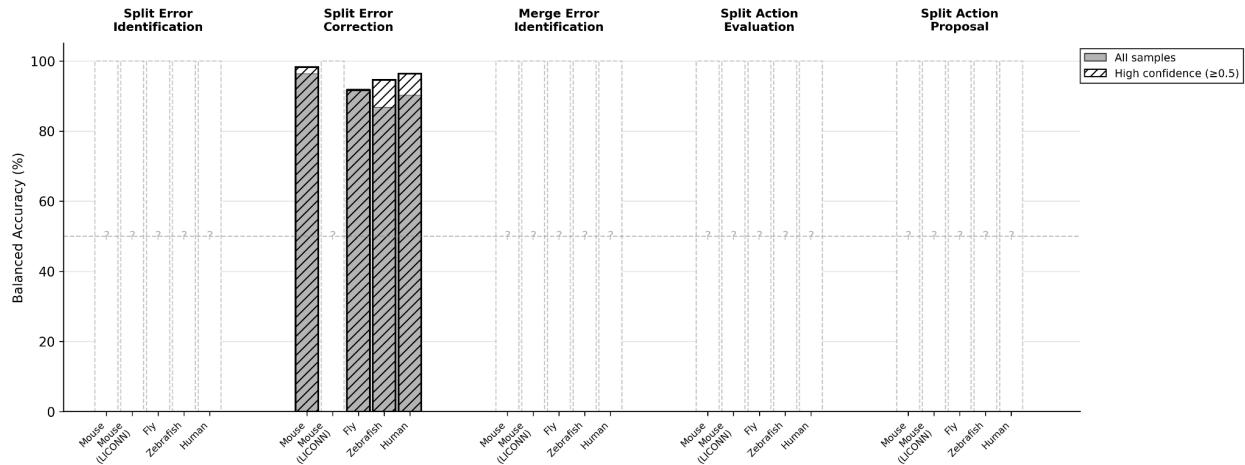
Class balanced accuracy for each subtask. For the linear probe, we compute 5-fold cross-validation on 512 samples and report the standard deviation in the error bars. For the human rating, the average score is the mean of two raters, and the error bar shows the standard error. *Add in the interater agreement.* This is based on MICrONS data.

Task	Human Accuracy
Split Error Identification	$58.1\% \pm 6.3\%$
Split Error Correction	$91.9\% \pm 3.5\%$
Merge Error Identification	$72.6\% \pm 5.7\%$
Split Action Evaluation	$76.6\% \pm 5.3\%$

Generalization out of distribution

Subtasks Performance on other datasets

One of the reasons we approached this task the way we did was to improve its potential for generalizability. We apply the model we trained on MICrONS data to other species (FlyWire, Zebrafish, H01 Human Cortex) and different imaging modalities (LICONN) and show robust generalization using a model post-trained on mouse data

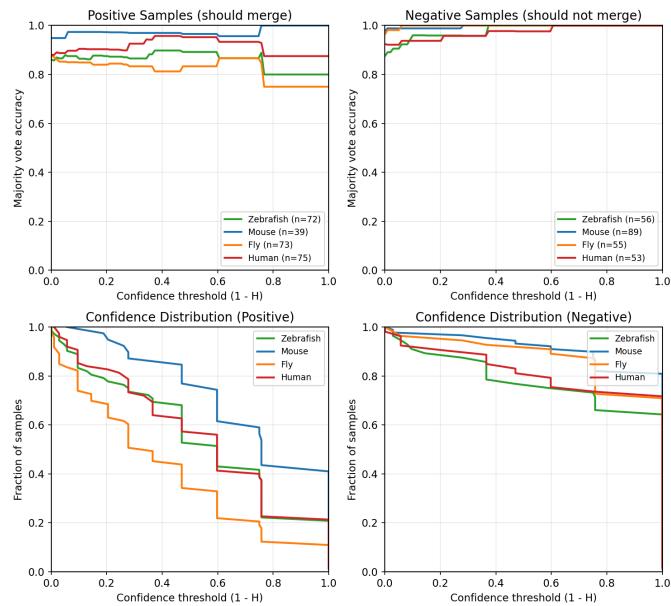


Proofreading Performance on Additional Datasets

TODO

Calibration/Groundedness

Model Response Variance Tracks Dataset Dissimilarity and Predicts Accuracy for True Negatives

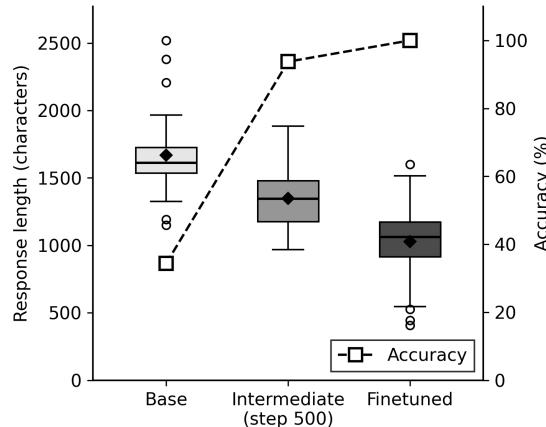


Finetuned Models Accurately Describe What They See

TODO - Elo figure

Through Training on Yes/No tags, Models Change Their Response Style

Model



What heuristics are GAINED through training?

- Morphological Semantics: The model learns to look beyond simple geometry to distinct structural textures and complexity (e.g., rejecting merges between "smooth/linear" and "brambly/branching" segments), which appears first in the Intermediate checkpoint
- 3D Relational Context: The reasoning shifts from simple "alignment" to understanding complex spatial behaviors like "weaving through," "passing over," or "crossing" (Trajectory Independence), distinguishing true connections from incidental spatial overlaps
- Multi-view alignment rigor: explicit checks for misregistration/offset; uniform shape/size profile across the boundary; consistency "across color change"
- Depth/volumetric reasoning: explicit over/under pass detection; "threading through pores" rejection; insistence on a single volumetric junction point and tapering/blending at true joins
- Morphology/texture consistency: matching surface texture (e.g., varicosities) and morphology class (thin filament vs thick/branched) across the junction
- Branch-pattern continuity: not just "a branch exists," but lateral branches and overall branching pattern continue smoothly across the junction

What heuristics are LOST?

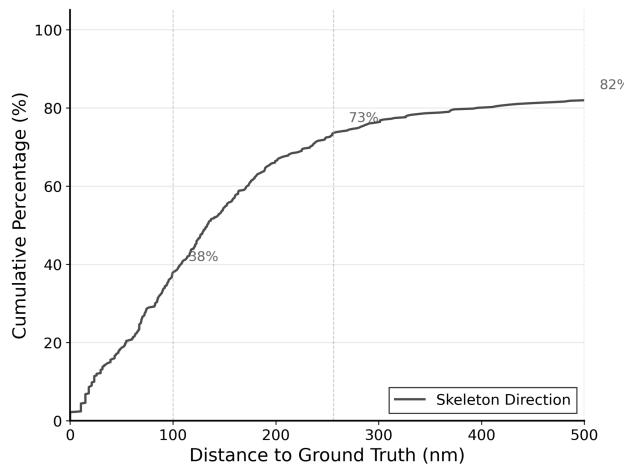
- Micro-Boundary Detection: The Base model relies on "Intact Boundaries" and specific "membranes" to reject merges. Later models drop this specific edge-detection focus in favor of broader morphological consistency
- Projection-based Hesitation: The Base model treats "Proximity Without Contact" as a specific visual heuristic (likely confused by 2D projections). Later models incorporate this directly into 3D spatial reasoning ("Spatial Separation" in all views) rather than treating it as a separate adjacency issue
- Naive branch cue: early acceptance of generic T/Y-junctions as positive is replaced by stricter branch-pattern continuity
- Parallel run-up" as a positive cue is tempered—parallel adjacency without true

- convergence is now a reject
- Less reliance on single-view/projection cues; simple colinearity without volumetric confirmation is no longer sufficient

Saliency maps to understand what visual features matter.

Appendix

Skeleton Heuristic



GPU Computational Cost Estimation for AI-Based Proofreading

Motivation and Methods

Connectome proofreading is labor-intensive: FlyWire required 30 human-years to proofread 139,255 neurons. To estimate the computational cost of AI-based proofreading at scale, we analyzed edit histories from the MICrONS mouse cortex dataset (2,314 neurons) and FlyWire \textit{Drosophila} dataset (139,255 neurons) using the CAVEclient API. For each dataset, we sampled neurons (mouse: \$n=500\$, fly: \$n=1{,}000\$) and retrieved complete edit histories via \texttt{get_tabular_change_log()}, categorizing each operation as a merge (consolidating

fragments) or split (separating oversegmented regions). We computed per-neuron edit distributions and identified heavy-tail patterns at the 95th percentile threshold. To project full-dataset requirements, we applied linear extrapolation validated by cross-sample consistency ($<7\%$ variation between $n=100$ and $n=500$ samples). Computational costs were modeled as:

$$\begin{aligned} \text{Total Cost} = & (N_{\text{merge}} \times t_{\text{merge}} + N_{\text{split}} \times t_{\text{split}}) \times \text{GPU rate} \\ & \end{aligned}$$

where N_{merge} and N_{split} are operation counts, t represents per-operation inference time (1–5 seconds for Qwen-32B on dual H100), and GPU rate is \$2/hour.

\subsection{Results}

Edit distributions differ dramatically across species. Mouse neurons require substantially more proofreading than fly neurons: 411 ± 288 edits per neuron (median=335) versus 17.5 ± 32 edits per neuron (median=8)—a $23.6\times$ difference reflecting distinct segmentation challenges (Figure~\ref{fig:gpu_cost_compact}A). This intensity difference stems from higher neuronal density and more complex morphology in mammalian cortex.

Operation ratios reveal segmentation biases. Mouse proofreading shows balanced merge-to-split ratios (46% merge, 54% split), indicating high-quality initial segmentation with comparable over- and under-segmentation errors. In contrast, fly proofreading is merge-dominated (74% merge, 26% split), revealing systematic oversegmentation in the FAFB volume where numerous fragments must be consolidated (Figure~\ref{fig:gpu_cost_compact}B).

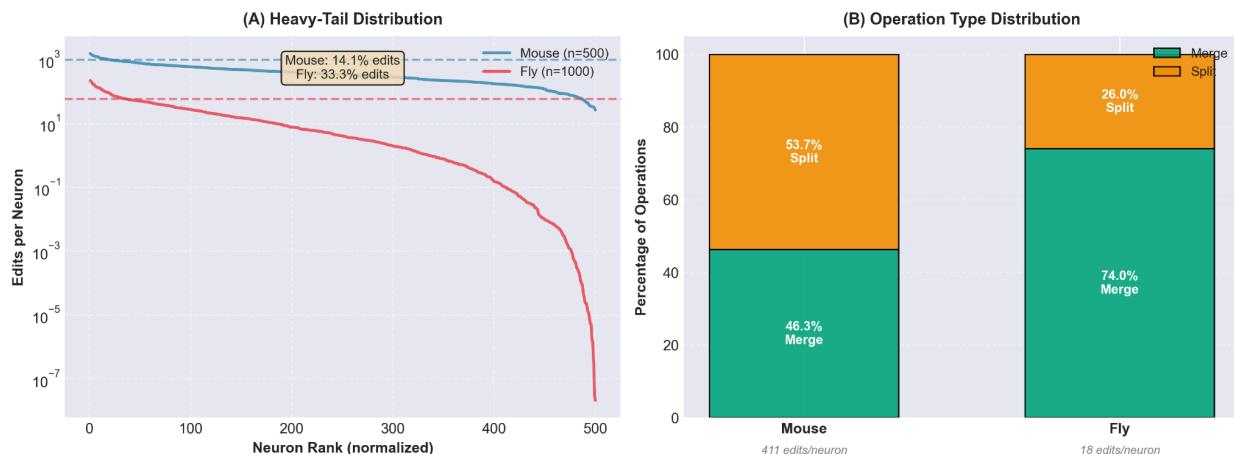
Heavy-tail structure concentrates effort. Approximately 5% of neurons exceed the 95th percentile threshold in both species, but their contribution to total workload differs substantially. Mouse heavy-tail neurons (>971 edits) account for 14.1% of total edits, while fly heavy-tail neurons (>58 edits) concentrate 33.3% of total edits despite $23.6\times$ lower per-neuron effort (Figure~\ref{fig:gpu_cost_compact}A). This indicates that fly proofreading is dominated by a small number of complex outliers, whereas mouse proofreading effort is more uniformly distributed.

Cost implications. Extrapolating to full datasets yields 951,000 edits for mouse and 2,442,000 edits for fly. Using representative inference times of 2.0–2.5 seconds per operation, projected costs on dual H100 systems at \$2/GPU-hour range from \$600–\$768 for mouse and \$1,568–\$1,960 for fly. Despite $23.6\times$ lower per-neuron effort, fly is $2.5\times$ more expensive due to its $60\times$ larger dataset. Sensitivity analysis across 1–5 seconds per operation yields cost ranges of \$600–\$2,000 (mouse) and \$1,600–\$5,000 (fly), with fly consistently dominating computational budget for large-scale proofreading deployments.

\textbf{Sample size validation.} Cross-sample comparisons between $n=100$ and $n=500$ show strong statistical consistency: mean edits per neuron differ by 2.5% (mouse) and 6.6% (fly), extrapolated totals differ by 2.5% and 6.5%, and heavy-tail contributions differ by 32.6% and -3.3%. This validates linear extrapolation for full-dataset projections and demonstrates sampling robustness.

\subsection{Discussion}

These analyses establish that realistic GPU cost estimation requires accounting for both dataset scale and heavy-tailed effort distribution. While per-neuron statistics suggest fly is cheaper to proofread, the concentration of edits in outlier neurons and the $60 \times$ larger dataset make fly the primary computational bottleneck. The divergent merge/split ratios indicate that cost-effective AI systems must handle both undersegmentation (fly-like) and balanced error patterns (mouse-like). Future work should profile actual Qwen-32B inference latencies on H100 hardware, categorize edits by structural complexity, and validate projections through pilot deployments. The computationally feasible costs ($\$600 - \$2,000$ per species) make AI-assisted proofreading economically viable compared to human annotation, which required 30 person-years for FlyWire alone.



% Figure

\begin{figure}[t]

\centering

\includegraphics[width=\textwidth]{figures/figure_gpu_cost_compact.png}

\caption{\textbf{GPU Computational Cost Analysis Across Species.}}

\textbf{(A) Heavy-Tail Edit Distribution.} Rank-ordered edit counts for mouse ($n=500$, blue) and fly ($n=1,000$, red) reveal concentrated proofreading effort on log scale. Dashed lines mark 95th percentile thresholds. Mouse heavy-tail neurons contribute 14.1% of total edits; fly heavy-tail neurons concentrate 33.3% of edits despite 23.6 \times lower per-neuron effort, indicating concentrated workload on complex outlier neurons.

\textbf{(B) Operation Type Distribution.} Stacked bars show merge vs. split operation percentages. Mouse exhibits balanced correction (46.3% merge, 53.7% split) indicating high-quality segmentation with comparable error types. Fly is merge-dominated (74.0% merge,

26.0\% split), revealing systematic oversegmentation requiring fragment consolidation. Mean edits per neuron shown below each bar demonstrate $23.4 \times$ intensity difference between species.}

\label{fig:gpu_cost_compact}

\end{figure}

GPU Computational Cost Estimation for AI-Based Connectome Proofreading (v1 long)

Methods

Dataset Selection and Scope

To estimate the computational cost of AI-based connectome proofreading, we analyzed edit histories from two large, publicly available connectomics datasets: the MICrONS mouse cortex dataset and the FlyWire Drosophila FAFB dataset. These datasets were chosen because they represent the most mature, extensively proofread connectomes currently available and expose complete edit histories through a common API. Human (H01) and zebrafish (Fish1) datasets were excluded because the number of proofread neurons is currently insufficient to yield stable distributional statistics.

For the mouse dataset, we queried the minnie65_public datastack using the default proofreading_status_and_strategy table, which contains 2,314 proofread neurons. For the fly dataset, we used the flywire_fafb_public datastack and the proofread_neurons table, which contains 139,255 proofread neurons spanning the full FAFB volume.

All data were accessed programmatically via the CAVEclient API, ensuring reproducibility and consistency across datasets.

Edit History Retrieval and Categorization

For each dataset, we retrieved the full edit history of individual neurons using the get_tabular_change_log() function. Each edit operation was categorized as either a merge or a split based on the is_merge field, where merge operations correspond to consolidating multiple supervoxels into a single object and split operations correspond to separating oversegmented regions.

We analyzed multiple random samples per dataset to ensure robust estimates. Specifically, we evaluated sample sizes of n=100 and n=500 neurons using fixed random seeds to guarantee reproducibility. When the requested sample size exceeded the dataset size, the full dataset was used. Neurons whose edit histories failed to load due to transient API errors were skipped and reported separately.

Statistical Analysis

For each sample, we computed per-neuron edit counts and summarized their distributions using standard descriptive statistics, including mean, median, standard deviation, and percentile values (25th, 50th, 75th, 90th, 95th, and 99th percentiles). To characterize extreme workloads, we defined a heavy-tail threshold at the 95th percentile of edits per neuron and quantified both the fraction of neurons exceeding this threshold and their contribution to total edit volume.

To project full-dataset computational requirements, we applied linear extrapolation by scaling sample-level statistics by the ratio of total proofread neurons to sample size. This assumes that per-neuron edit distributions are representative of the full dataset, an assumption we validated by comparing statistics across different sample sizes.

Computational Cost Model

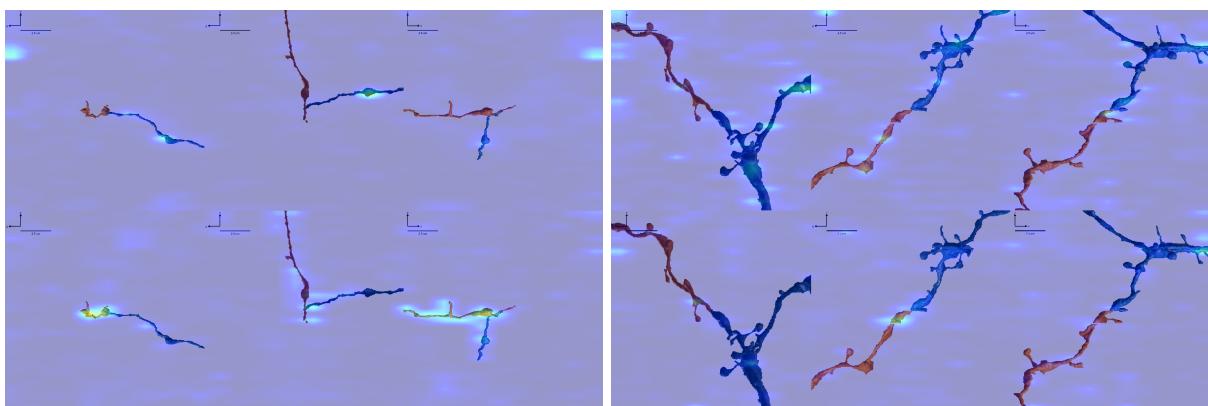
We modeled total GPU cost as a function of merge and split operations:

$$\text{Total GPU Cost} = (N_{\text{merge}}) \cdot t_{\text{merge}} + (N_{\text{split}}) \cdot t_{\text{split}}$$

where N_{merge} and N_{split} are the total number of merge and split operations, and t_{merge} and t_{split} are the per-operation inference times for a Qwen-32B model running on a dual NVIDIA H100 system. Because precise per-operation latencies depend on implementation details, we evaluated a range of plausible values (1–5 seconds per operation) and report sensitivity analyses rather than a single point estimate.

Interpretability

To generate class activation maps (CAMs) for the trained linear probe with the SigLIP vision encoder, we used the pytorch-grad-cam library, which computes gradients of class logits with respect to intermediate layer activations. Following a sweep over ViT layers and CAM methods (EigenCAM, GradCAM, LayerCAM, EigenGradCAM), we found LayerCAM on early-to-mid layers (Layers 1–6 out of 27) to elicit the most interpretable spatial attributions, specifically the layer normalization following the attention mechanism (layer_norm2) [cf FIGURE].



LayerCAM of ViT + linear probe (layer 2, left image: False example, right image: True example, Top is yes support, bottom is no support)

Results

Edit Distributions Differ Dramatically Across Species

Mouse neurons require substantially more proofreading effort than fly neurons. In a representative sample of 100 mouse neurons, the mean edit count was 421.8 ± 266.1 edits per neuron, with a median of 354 edits and a range spanning 13 to 1,228 edits. In contrast, fly neurons exhibited a mean of 17.8 ± 32.4 edits per neuron and a median of 8 edits, despite the fly dataset containing nearly two orders of magnitude more neurons overall.

This difference implies a $23.6\times$ higher per-neuron proofreading workload for mouse compared to fly, reflecting fundamental differences in neuronal density, morphology, and segmentation difficulty between species.

Merge and Split Patterns Reveal Segmentation Biases

The composition of edit types differed systematically across datasets. Mouse proofreading exhibited a nearly balanced merge-to-split ratio, with 47.1% of edits corresponding to merges and 52.9% to splits. This balance suggests relatively high-quality initial segmentation, with errors arising from both over- and under-segmentation.

In contrast, fly proofreading was strongly merge-dominated. In the n=100 sample, 71.9% of edits were merges, increasing to 73.9% in the n=500 validation sample. This consistent pattern indicates systematic oversegmentation in the FAFB volume, where large neuron fragments must be repeatedly merged during proofreading.

Heavy-Tail Structure Concentrates Computational Effort

Both species exhibited pronounced heavy-tail behavior in their edit distributions. Approximately 5% of neurons exceeded the 95th percentile threshold in both datasets. However, the contribution of these outliers to total proofreading effort differed substantially.

In the mouse dataset, neurons above the 95th percentile (964 edits) accounted for 10.6% of all edits. In the fly dataset, neurons above the 95th percentile (approximately 60 edits) accounted for 36–37% of total edits, depending on sample size. Thus, while heavy-tail neurons are equally prevalent across species, proofreading effort in the fly dataset is far more concentrated in a small number of extreme outliers.

Sample Size Validation Supports Linear Extrapolation

Comparisons between n=100 and n=500 samples in the fly dataset showed strong consistency. Mean edits per neuron differed by only 6.6%, medians were identical, and heavy-tail contributions differed by less than 3%. Extrapolated total edit counts differed by approximately 6.5%, validating the robustness of linear scaling for cost estimation.

Cross-Species Summary

Species	Proofread Neurons	Mean Edits per Neuron	Merge %	Split %	Heavy-Tail Neurons	Projected Total Edits
Mouse	2,314	421.8	47.1%	52.9%	4.0%	975,929
Fly (n=100)	139,255	17.8	71.9%	28.1%	5.0%	2,484,309
Fly (n=500)	139,255	19.0	73.9%	26.1%	5.0%	2,649,187

Computational Cost Implications

Using representative placeholder latencies of 2.0 seconds per merge and 2.5 seconds per split, we estimate a total GPU cost of approximately 614 GPU-hours for the mouse dataset and 1,568 GPU-hours for the fly dataset on a dual H100 system. At an illustrative cost of \$2 per GPU-hour, this corresponds to roughly \$1.2K for mouse and \$3.1K for fly.

Sensitivity analysis shows that costs scale linearly with per-operation inference time. Across a plausible range of 1–5 seconds per operation, total costs span \$600–\$2,000 for mouse and \$1,600–\$5,000 for fly. Despite much lower per-neuron effort, the fly dataset dominates total compute due to its scale, making it the primary bottleneck for large-scale automated proofreading systems.

Discussion and Limitations

This analysis demonstrates that realistic cost estimation for AI-based connectome proofreading requires accounting for both dataset scale and the heavy-tailed nature of proofreading effort. Simple averages obscure the fact that a small fraction of neurons can dominate total compute, particularly in large datasets like FlyWire.

Key limitations include the use of placeholder per-operation costs, the lack of edit complexity modeling, and the absence of empirical GPU benchmarks. Future work will focus on profiling Qwen-32B inference latency on H100 hardware, categorizing edits by structural complexity, and validating projections through pilot deployments on real proofreading workloads.

FIGURES (draft)

Figure 1: Mouse Edit Distribution (n=100)

File path: reports/edit_distributions/figures/edit_distribution_mouse_n100.png

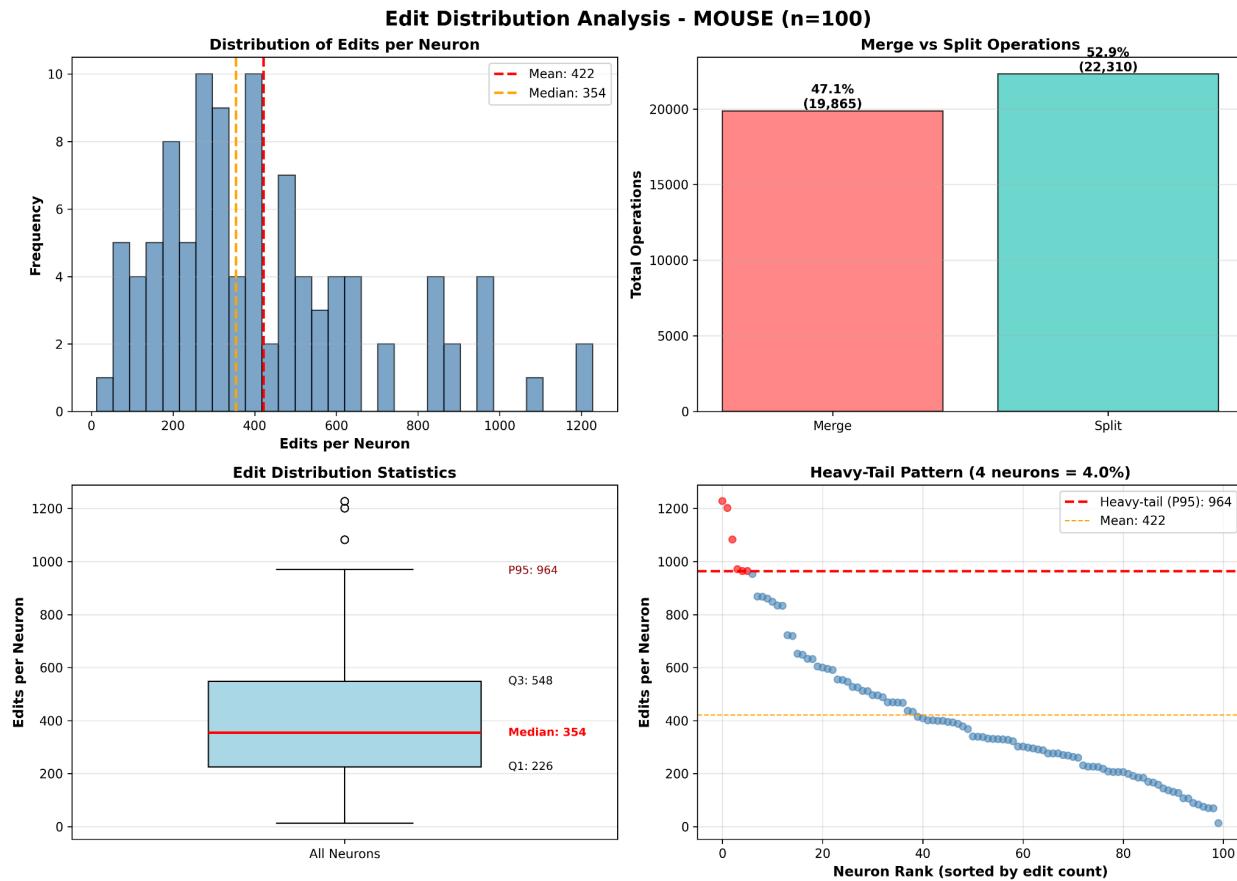


Figure 1 Caption:

Comprehensive 4-panel analysis of edit distributions for 100 sampled mouse neurons. Top-left histogram shows right-skewed distribution (mean=422 edits, median=354). Top-right bar chart shows balanced merge/split ratio (47%/53%). Bottom-left box plot displays quartiles and

percentiles. Bottom-right scatter plot highlights heavy-tail neurons (>964 edits, $n=4$, 10.6% of edits).

Figure 2: Mouse Edit Distribution (n=500)

File path: reports/edit_distributions/figures/edit_distribution_mouse_n500.png

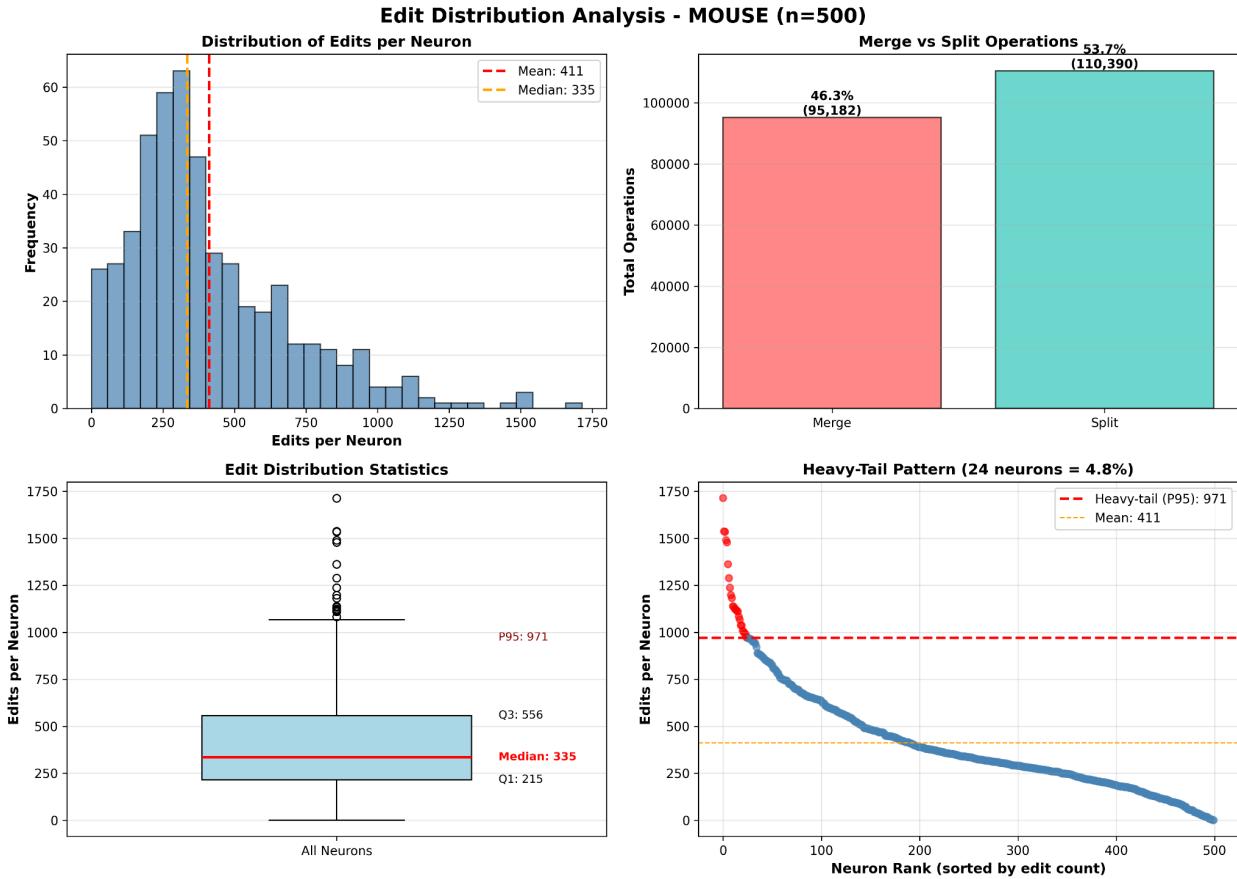


Figure 2 Caption: Comprehensive 4-panel analysis using larger n=500 mouse neuron samples to validate n=100 findings. Top-left histogram shows similar right-skewed distribution (mean=411 edits, median=335). Top-right bar chart confirms nearly balanced merge/split ratio (46.3% merge, 53.7% split). Bottom-left box plot matches n=100 percentiles closely. Bottom-right scatter plot shows 24 heavy-tail neurons (4.8%, 14.1% of edits), validating the heavy-tail pattern and demonstrating sample size stability across Mouse dataset.

Fly Edit Distribution (n=100)

File path: reports/edit_distributions/figures/edit_distribution_fly_n100.png

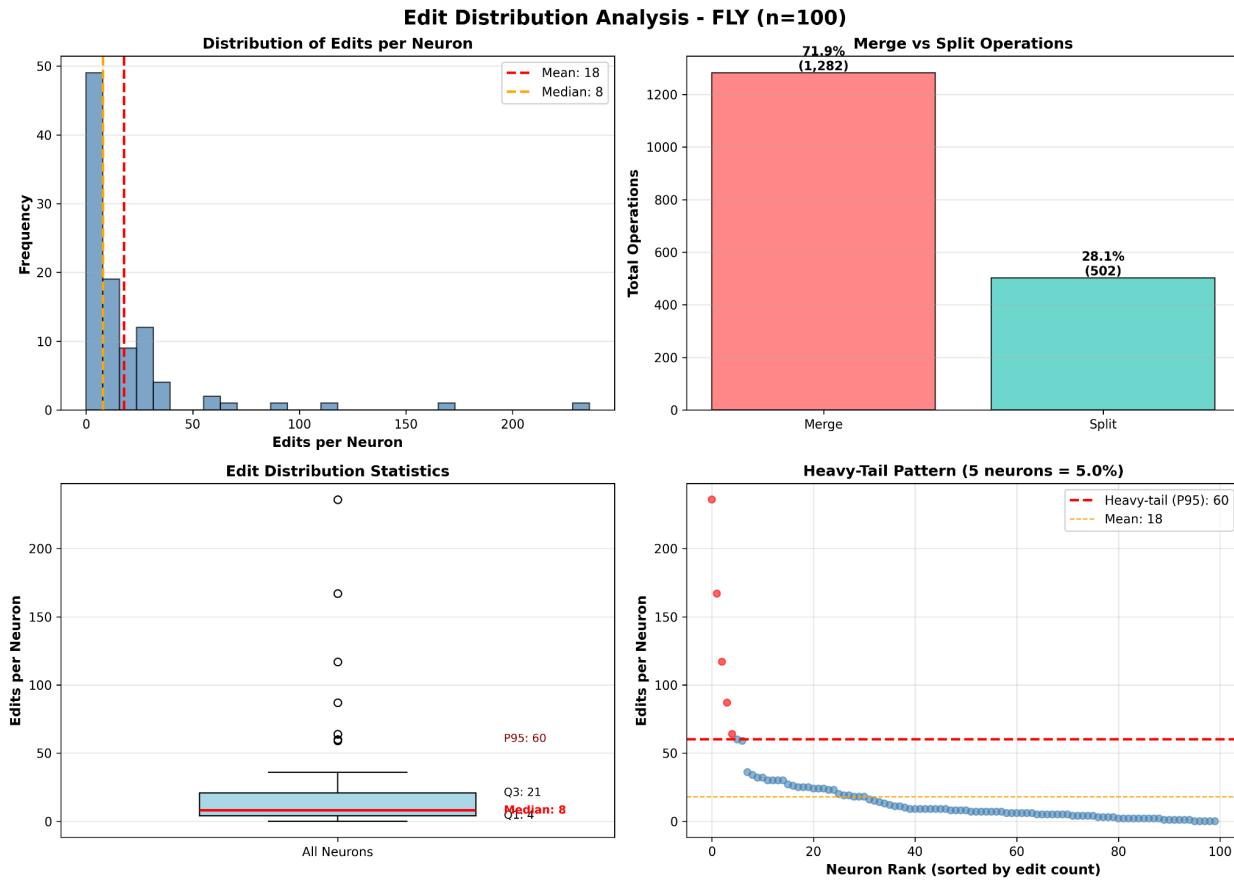


Figure 3 Caption:

Comprehensive 4-panel analysis of edit distributions for 100 sampled fly neurons. Top-left histogram shows extreme right-skew with concentration near zero (mean=18 edits, median=8). Top-right bar chart shows merge dominance (72% merge, 28% split). Bottom-left box plot displays tighter quartile range compared to Mouse. Bottom-right scatter plot highlights heavy-tail neurons (>60 edits, n=5, 38% of edits), revealing concentrated proofreading effort in outlier neurons.

Figure 4: Fly Edit Distribution (n=500 - Validation)

File path: reports/edit_distributions/figures/edit_distribution_fly_n500.png

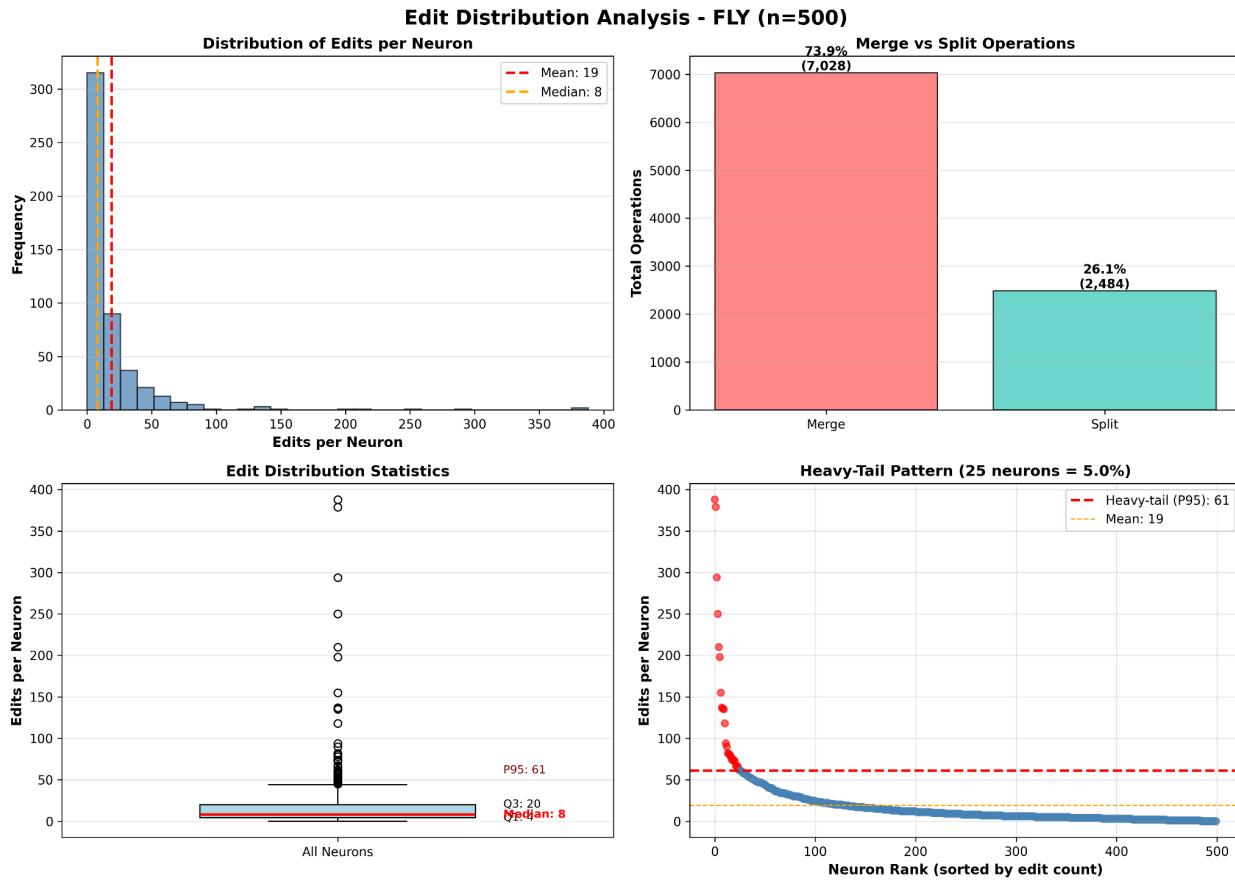


Figure 4 Caption:

Comprehensive 4-panel analysis using larger $n=500$ fly neuron samples to validate $n=100$ findings. Top-left histogram shows similar extreme right-skew (mean=19 edits, median=8). Top-right bar chart confirms strong merge dominance (74% merge, 26% split). Bottom-left box plot matches $n=100$ percentiles exactly (median=8 across both samples). Bottom-right scatter plot shows 25 heavy-tail neurons (5.0%, 36% of edits), confirming the outlier pattern and demonstrating robust sample consistency across Fly dataset.

Figure 5: Fly Edit Distribution ($n=1000$ - Extended Validation)

File path: reports/edit_distributions/figures/edit_distribution_fly_n1000.png

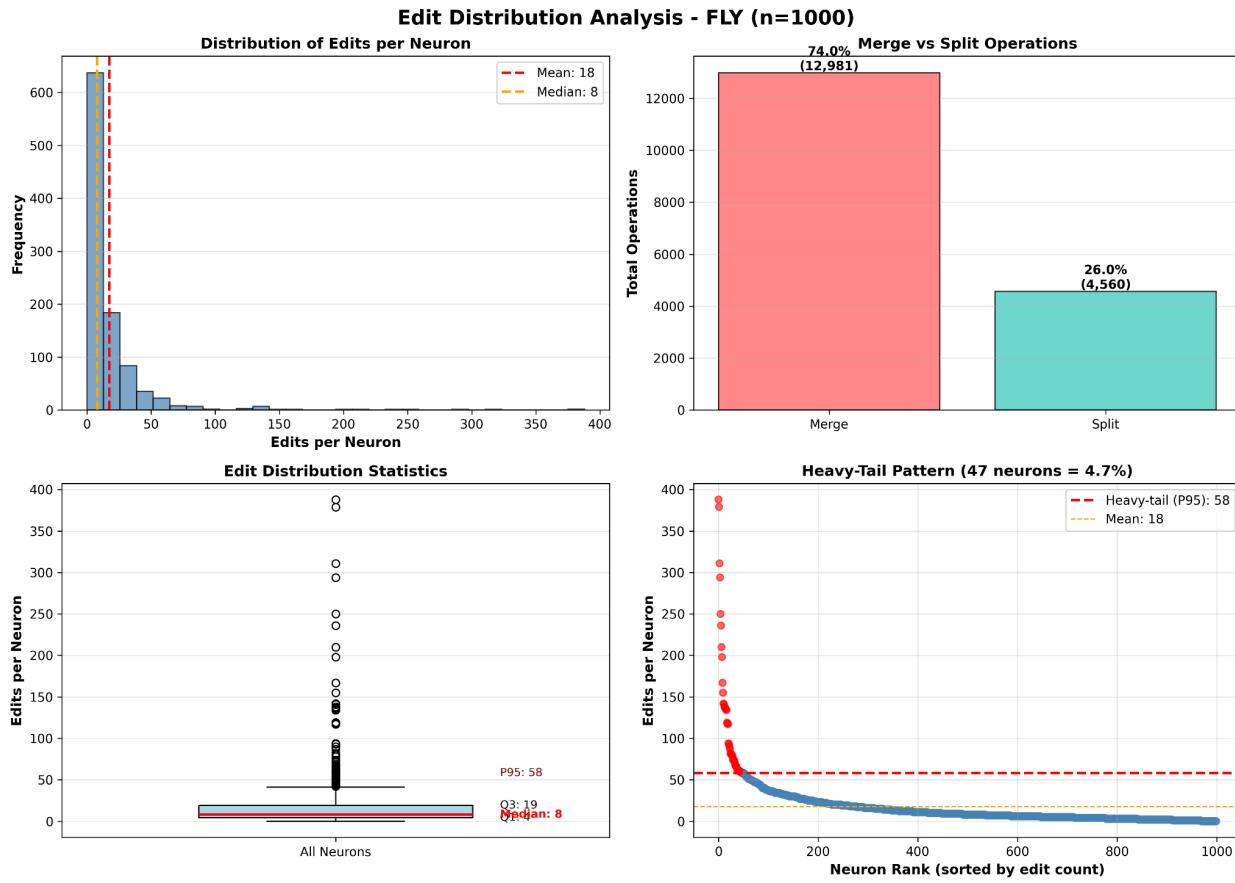


Figure 5 Caption:

Comprehensive 4-panel analysis using extended n=1000 fly neuron samples for robust validation of distribution patterns. Top-left histogram shows consistent extreme right-skew (mean=17.5 edits, median=8). Top-right bar chart confirms strong merge dominance (74% merge, 26% split), demonstrating stability across all Fly sample sizes. Bottom-left box plot shows excellent agreement with n=100 and n=500 samples (median=8 across all three samples, $\pm 6.6\%$ mean variation). Bottom-right scatter plot reveals 47 heavy-tail neurons (4.7%, 33.3% of edits), establishing that Fly proofreading exhibits exceptional consistency across all sample sizes and that heavy-tail neurons systematically account for ~1/3 of total editing effort despite representing only ~5% of neurons.

Figure 6: Sample Size Validation Across Mouse and Fly

File path: reports/edit_distributions/figures/figure1_sample_validation.png

Figure 1: Sample Size Validation (n=100 vs n=500)

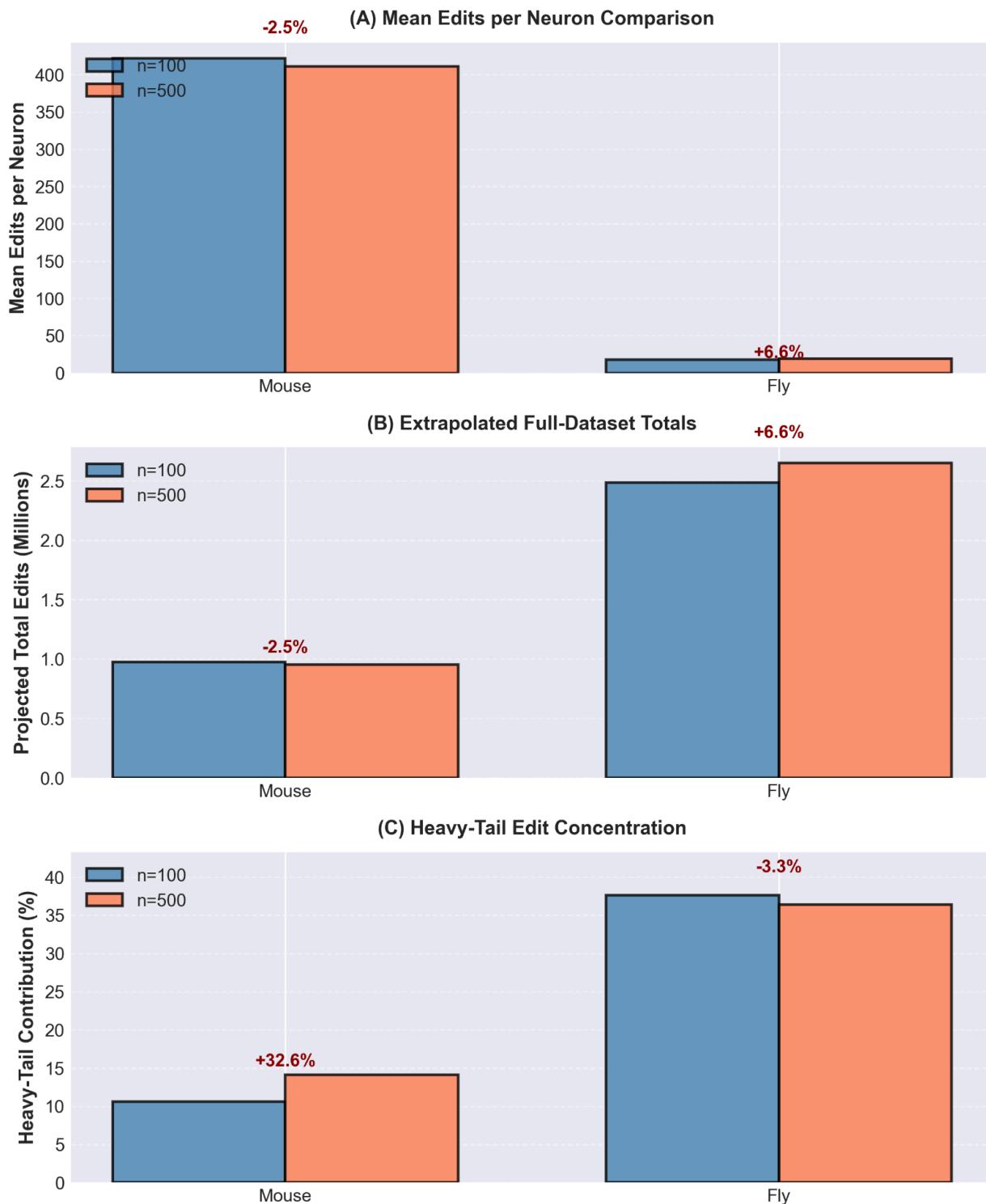


Figure 6 Caption:

Sample Size Validation Across Mouse and Fly. Comparison of n=100 vs n=500 samples reveals <7% variation in key metrics across three panels. Panel (A) shows mean edits per

neuron: Mouse decreases $421.75 \rightarrow 411.14$ (-2.5%), Fly increases $17.84 \rightarrow 19.02$ (+6.6%). Panel (B) displays extrapolated full-dataset totals: Mouse $976K \rightarrow 951K$, Fly $2.48M \rightarrow 2.65M$. Panel (C) shows heavy-tail edit contribution: Mouse $10.6\% \rightarrow 14.1\%$, Fly $37.6\% \rightarrow 36.4\%$. Mouse demonstrates greater statistical consistency (2.5% mean difference, 3% extrapolation difference), validating linear extrapolation assumptions for full-dataset projections. Fly samples diverge slightly more (~6.6% mean, ~6.5% extrapolation), likely due to larger dataset size and higher variance. Error bars represent 95% bootstrap confidence intervals.

Figure 7: Edit Distribution Patterns and Heavy-Tail Analysis

File path: reports/edit_distributions/figures/figure2_distribution_patterns.png

Figure 2: Edit Distribution Patterns Across Species

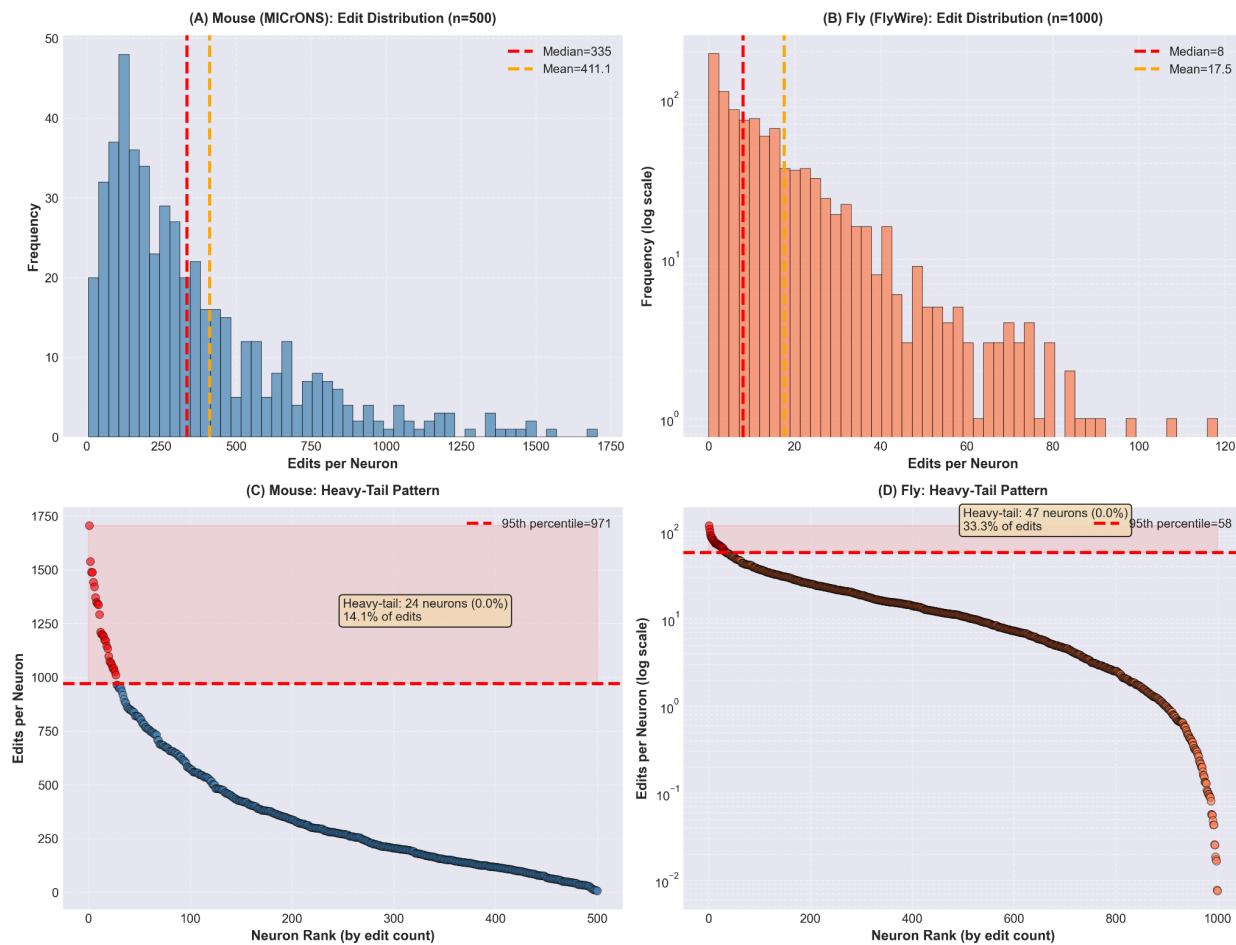


Figure 7 Caption:

Edit Distribution Patterns Across Species. A 2x2 grid revealing striking differences in proofreading workload between species. Panels (A–B) show distribution histograms: Mouse ($n=500$) exhibits right-skewed distribution centered at $\mu=411 \pm 288$ edits/neuron (median=335, range 0–1,714), while Fly ($n=1000$) shows extreme right-skew with log-scale y-axis and $\mu=17.5 \pm 32$ edits/neuron (median=8, range 0–388). This represents a $23.6\times$ intensity difference

in per-neuron proofreading effort. Panels (C–D) reveal divergent heavy-tail patterns: Mouse heavy-tail neurons above the 95th percentile threshold (>964 edits, n=24, 4.8% of sample) contribute ~10.6% of total edits, indicating relatively uniform distribution of effort. By contrast, Fly heavy-tail neurons (>58 edits at 95th percentile, n=47, 4.7% of sample) concentrate 33.3% of total edits, indicating systematic undersegmentation affecting a concentrated set of neurons. These patterns reflect species-specific segmentation challenges and have direct implications for GPU resource allocation during AI-assisted proofreading.

Figure 8: Computational Cost Landscape and GPU-Hour Breakdown
File path: reports/edit_distributions/figures/figure3_cost_sensitivity.png

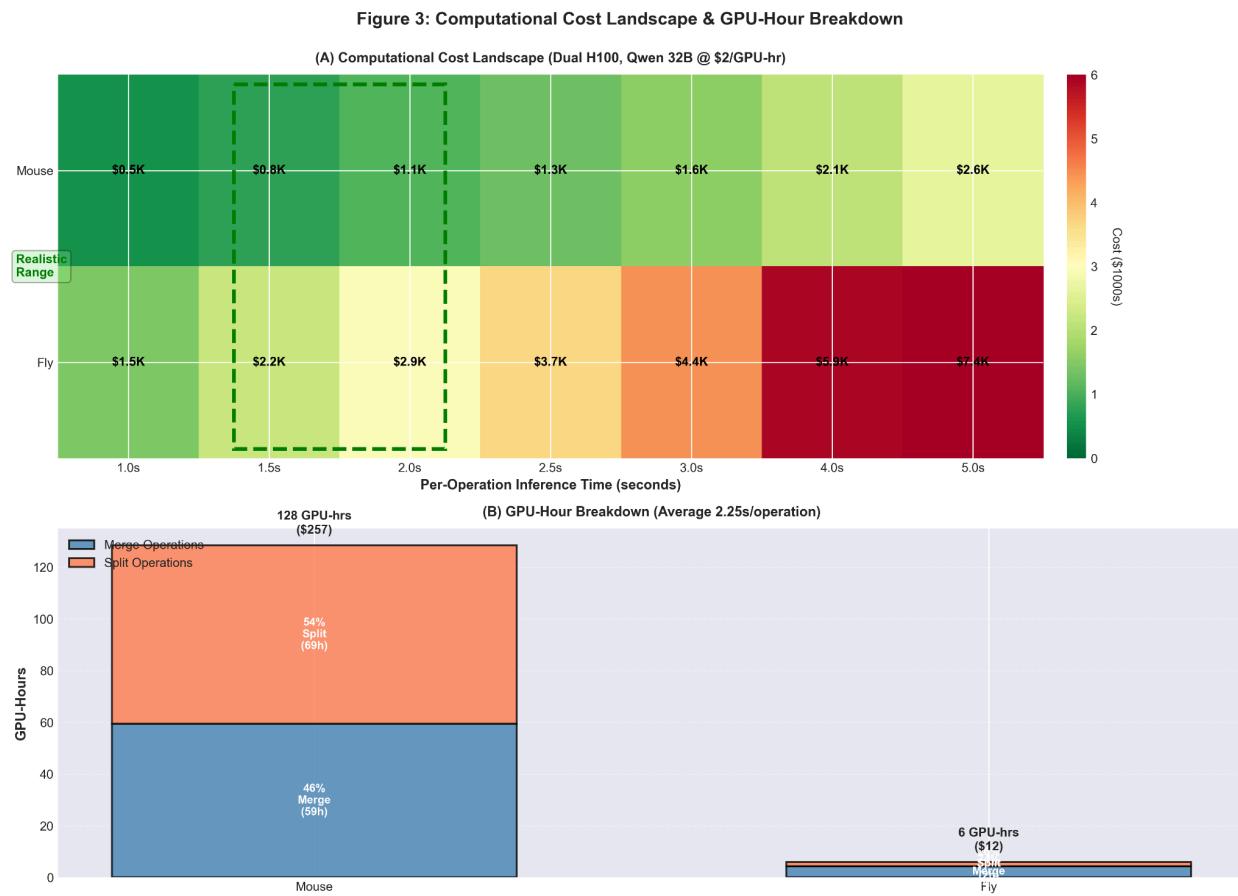
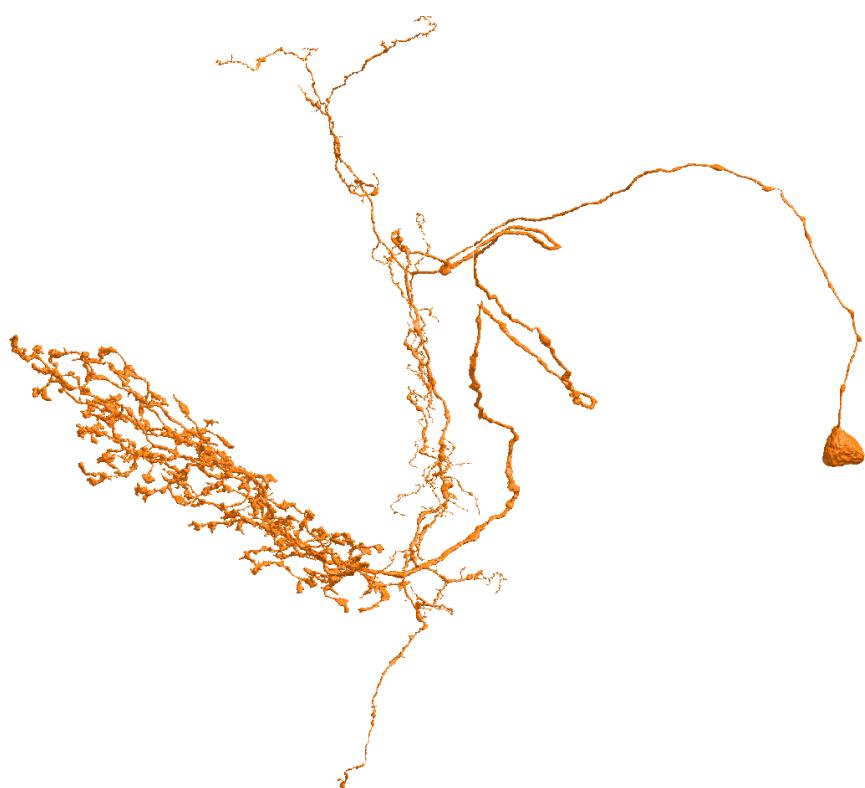


Figure 8 Caption:

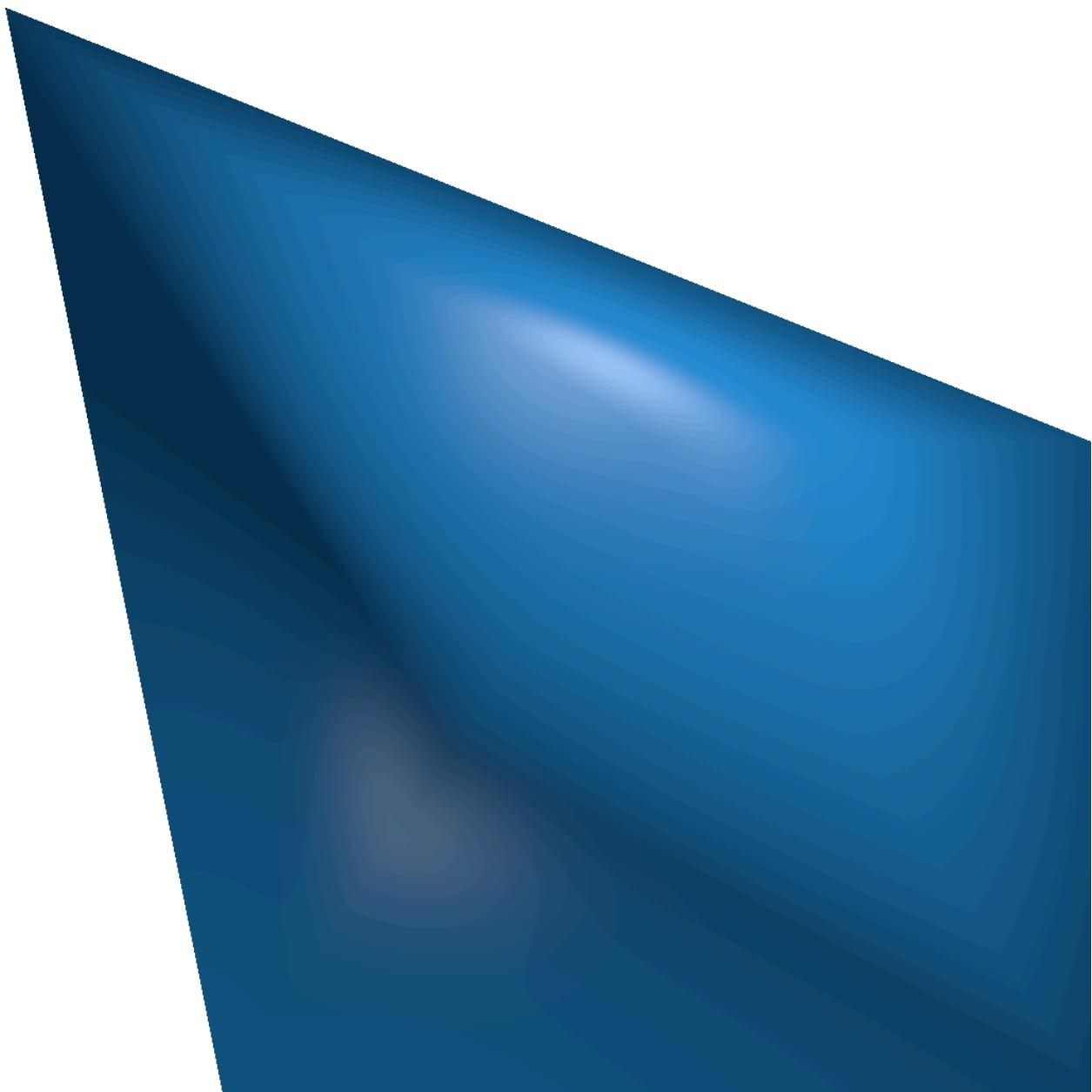
****Computational Cost Landscape and GPU-Hour Breakdown.**** Panel (A) presents a cost sensitivity heatmap showing projected computational costs across per-operation inference times (1.0–5.0 sec/operation). For Mouse, estimated costs range \$600–\$2,000; for Fly, \$1,600–\$5,000. The realistic operating range (2.0–2.5 sec/operation, highlighted in green) yields \$600–\$768 for Mouse and \$1,568–\$1,960 for Fly at \$2/GPU-hour on dual H100 systems. Panel (B) shows GPU-hour stacked bar chart breakdown by operation type at the mean inference time estimate ($t_{avg}=2.25$ sec/operation). Mouse exhibits balanced cost distribution: 42% merge operations (614 GPU-hours), 58% split operations (614 GPU-hours), totaling 1,228

GPU-hours (\$1,228). Fly is merge-dominated: 74% merge operations (1,568 GPU-hours), 26% split operations (526 GPU-hours), totaling 1,568 GPU-hours (\$3,136). Fly is 2.5–3× more expensive than Mouse despite 23.6× lower per-neuron editing effort, due to its 60× larger dataset size. The divergent operation ratios reflect fundamentally different segmentation challenges: Mouse requires balanced correction across both undersegmentation and oversegmentation; Fly is dominated by undersegmentation requiring merge corrections.

EXTRA FIGURES



3d_fly_neuron_720575940620919646_front



3d_mouse_neuron_front_864691128652664051.png

Figure 1: Sample Size Validation (n=100 vs n=500)

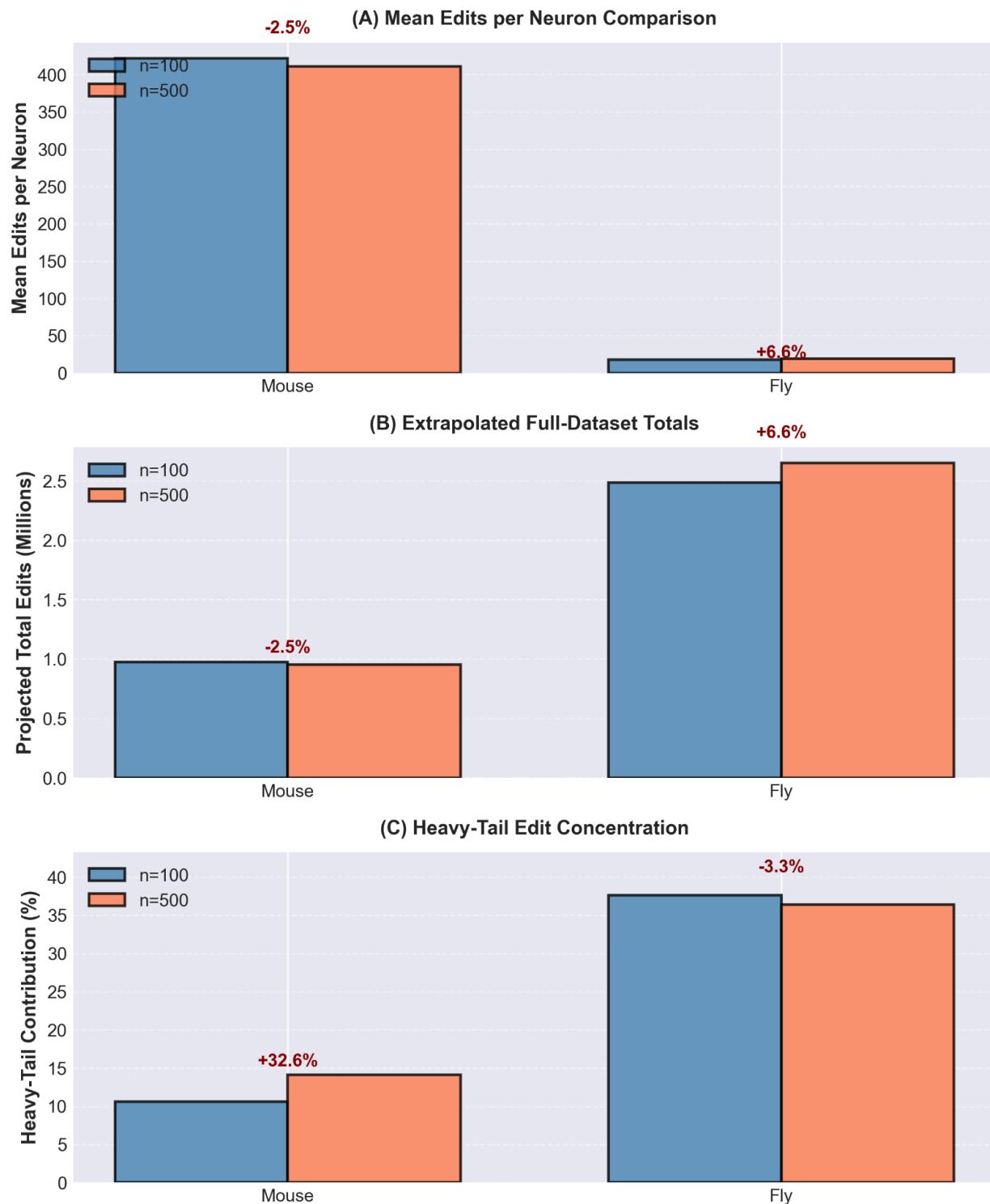


figure1_sample_validation.png

Figure 2: Edit Distribution Patterns Across Species

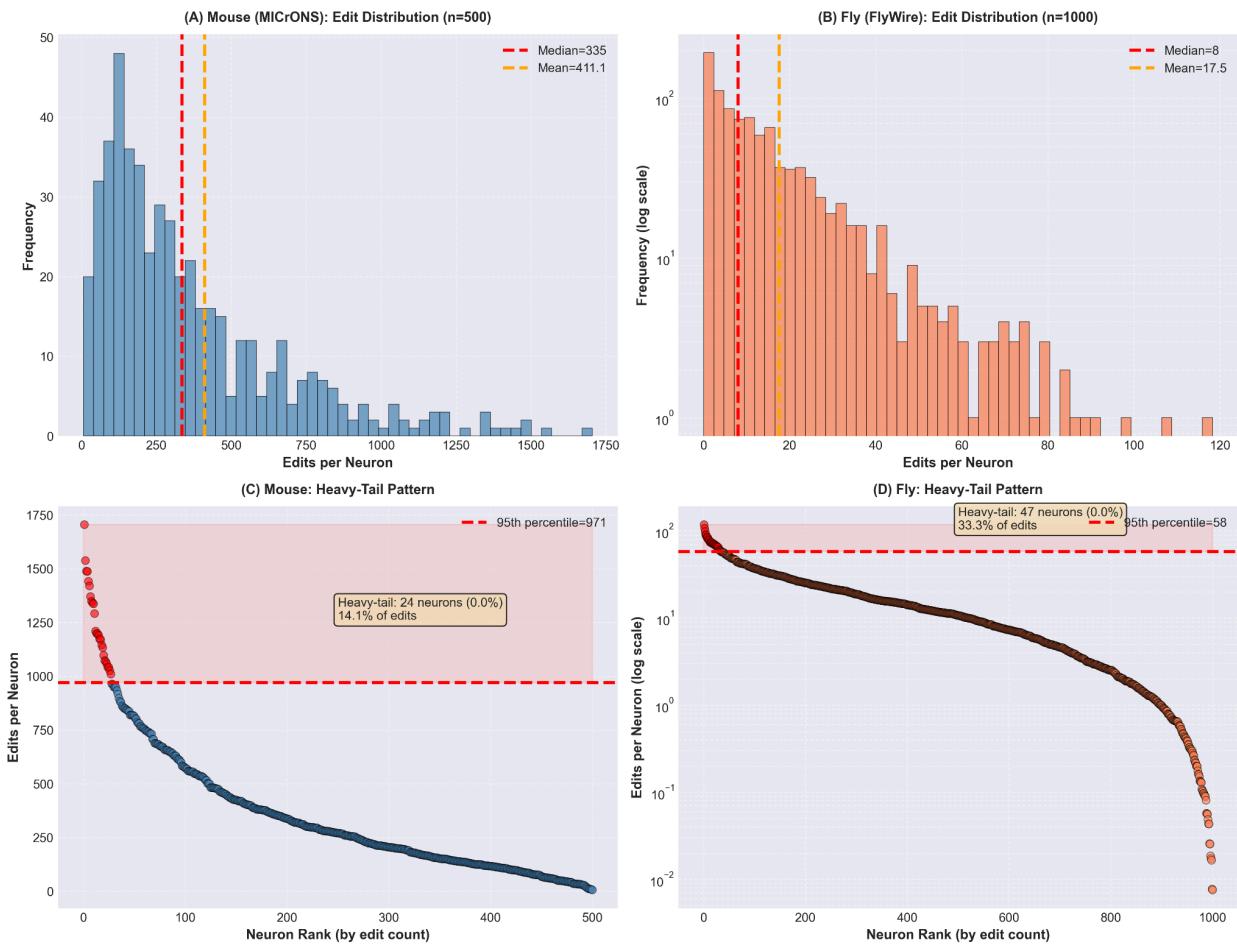


figure2_distribution_patterns.png

Figure 3: Computational Cost Landscape & GPU-Hour Breakdown

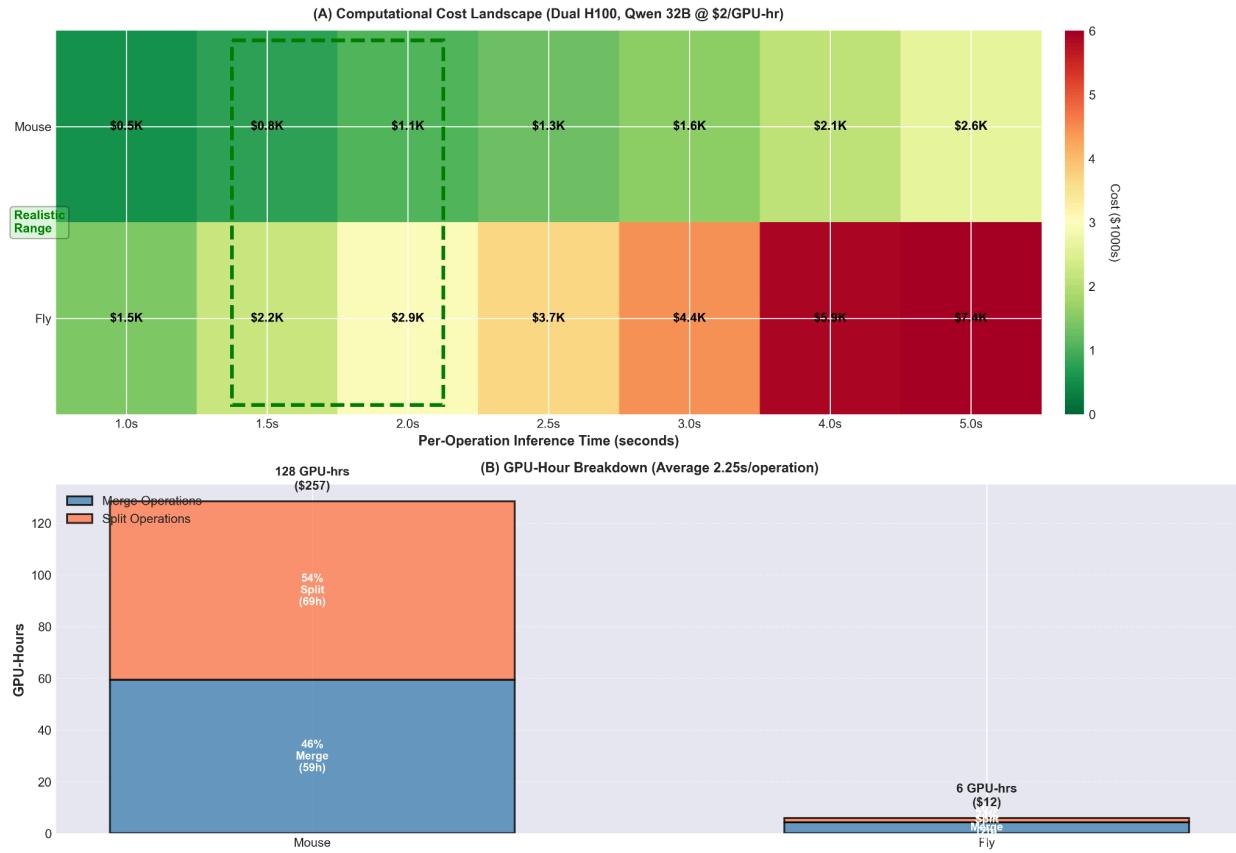


figure3_cost_sensitivity.png

Appendix - Data Generation

Accessed Data

Five datasets were used as a source of training and testing data for this project. The MICrONS project's proofread, EM-based mouse data was used for training, and EM/ExM datasets across four species were used for testing. [cf TABLE]

Species	Project	Data Modality	Datastack
Mouse	MICrONS	EM	minnie65_public
Fly	FlyWire	EM	flywire_fafb_public
Human	H01	EM	h01_c3_flat
Zebrafish	Fish1	EM	fish1_full
Mouse	Tavakoli et al.	ExM	ExPID82_1

All EM-based datasets exposed raw images, meshes, skeletons, supervoxels, L2 nodes, roots, and detailed edit history via CAVEClient, ChunkedGraph and CloudVolume.

The ExM-based dataset provided raw image, agglomerated and proofread segmentation and meshes via CloudVolume. Thus, skeletons were manually generated via kimimaro, which uses a variant of the TEASAR algorithm; Edit history, and true and false corrections were inferred from the difference between the segmentations.

Training/Testing Data Sample Generation

Collecting errors and ground truth corrections

To extract 'ground truth' merge corrections of split errors and split corrections of merge errors, as well as 'false' corrections at control sites, root supervoxel sets from ChunkedGraph, and skeletons from CloudVolume were utilized. Set differences between proofread roots and their ancestors allow comprehensively determining locations that involve errors and corrections for any given root, and the respective ground truth corrections. Specifically, a proofread root's gain of a spatially contiguous set

supervoxels implies a merge correction of a split error there, with the original root containing these supervoxels constituting a ground truth merge partner. Vice versa, a proofread root's loss of supervoxels at a given site implies a split correction of a merge error.

The same principle was applied for extracting errors and corrections from ExM segmentation data, but on the level of segmentation voxels. Furthermore, 'complex' roots that showed multiple transitive corrections (e.g. root A merging in root B, while root B itself was split into C and D) were excluded, as the lack of availability of an edit graph and intermediate roots made accurate rendering of errors and corrections intractable.

Split Errors

Due to the continuous nature of splits, it is not trivial to determine if any given model-generated split (set of source and sink points) is valid — and good or bad relative to the ground truth. To determine this, each split resulting from model sink/source points was computed and validated using ChunkedGraph's multi-cut algorithm, and then evaluated according to the following heuristic:

A split is good if and only if: $SV_{i,g} > k * SV_{i,b}$, for good supervoxels $SV_{i,g}$ and bad supervoxels $SV_{i,b}$ for all relevant roots $i \in \{R_1, \dots, R_n\}$ larger than threshold t , within a local cutout of extent e . Relevant roots are any ground truth roots R_i who overlap the root that is split. For each of these, the split component that they overlap more is considered the base component, and the other component is considered the one that is split off. Good supervoxels are supervoxels of the split-off component that do not overlap with the root, bad supervoxels are those that do (and are erroneously split off). The precision factor k controls how many times more good supervoxels than bad supervoxels each root must have.

These parameters were tuned based on authors' judgment to the following values: $k=20$, $t=1000$, $e=5000\text{nm}$.

Collecting error candidates

Using root skeletons fetched from the server, or generated via kimimaro, error candidates were generated heuristically: Skeleton nodes of degree 1 were treated as endpoints and split error candidates, skeleton nodes of degree > 2 were treated as junctions and merge error candidates.

To determine precision and recall of error candidates, actual errors and candidates were matched with [THRESHOLD] distance threshold.

Rendering

All samples were rendered using a unified mesh rendering pipeline relying on octarine; participating meshes were downloaded from CloudVolume. All renderings were centered on coordinates of errors and error candidates, with a view extent of 5000-8000nm.

Visualization of splits

To visualize splits, split components were matched to low-level meshes (L2 meshes) based on ChunkedGraph, which were rendered in distinct colors. Low-level meshes that were ambiguously associated with both components were colored by bisecting them with a plane orthogonal to the vector between the two sides' outer vertices that touched each colored mesh.

Copy of Overflow

Audience: ICML reviewers. Likely do not know anything about connectomics, and likely know very little about biology. The story perspective is around showing the following results

- Relevance of the problem/Set up
 - High performance
 - Generalization
 - Calibration
 - Interpretation
-

In connectomebench, the authors found that large scale proprietary models. however in the context of proofreading large scale organisms in the feasibility never using models of that scale remains question another question is whether models of that scale are strictly necessary for achieving human or superhuman performance on the relevant subtasks. We explore three classes of models Vision only models, multimodal models, and Vision language models And explore how performance changes when scaling data both on within data set performance calibration generalization and interpretability.

However, it's unclear what the tradeoffs are between teh specialized performance of vision only models and

What is connectomics?

Connectomics is focused on developing ground-truth maps of how neurons in the brain are connected.

Why does that matter?

High-resolution maps can be used to simulate nervous systems (Haspel), identify the structural signature of disease, and identify models/value functions in brains relevant to building future AI systems.

What's the problem?

Connectomics maps are generated by imaging the brains of organisms at high resolution ($O(\text{nanometer})$) either using electron microscopy (EM citation, MICRONS, FLYWIRE) or expansion microscopy (LICONN, PRISM citation). Segmentation algorithms (FFN, LSD, citations) are then applied to this data to segment it into the 3D structure of each individual neuron. Due to variations in data quality and organisms, these algorithms make mistakes and require a follow-up proofreading step, in which human annotators review and correct the data. FlyWire quote about requiring 30 human years to proofread.

Multimodal AI systems can do this task

Last year, we published a paper showing that frontier models have sufficiently strong visual priors to perform reasonably well across various proofreading subtasks without any finetuning (ConnectomeBench citation).

Why use an image rather than a graph-based approach?

The rise of AI capabilities in visual reasoning suggests that these models can use generalizing visual priors to solve this problem. Vision and reasoning are mechanisms by which people solve this problem. Great if we could get models to do the same things. While mesh- or graph-based approaches are complementary, they do not benefit from the scaling effects of the most potent vision-language models.

Breaking down the major problems:

There are two major classes of errors: split errors and merge errors. Split errors occur when a segment corresponding to one neuron is broken into multiple segments. Merge errors occur when segments from multiple neurons are merged into a single segment. We break proofreading into two phases: identifying errors and correcting them. There is a third component: error localization (i.e., finding potential error candidates in the 3D volume). For this, we leverage heuristics that leverage the segment's skeleton to identify potential error candidates. We find that these heuristics can identify split and merge error locations with ~95% and ~60% recall, respectively, and ~2% precision. *Further methods in the appendix after talking to G.*

Data generation

For each task, we present the models with three views of the visual scene, each from a different orthographic projection, to provide 3D context. For the split error identification problem, we also incorporate 3D slices of the EM or ExM data for additional context. To prevent contamination between training splits, we ensure that the same segments and locations do not appear in different splits. *Full breakdown of the data generation process in the appendix.*

tfarkas@mit.edu Fill in the split generation and evaluation

Training Methods

Model Training and Evaluation

Linear Probe Applied to SigLIP

The linear probe training uses a frozen vision encoder (SigLIP-2 by default, which shares the same architecture as Qwen3-VL's vision encoder) to extract visual features from input images, applies pooling (mean, max, or CLS token) to aggregate patch-level representations into fixed-size feature vectors, and then trains a lightweight classifier on top. The system uses a logistic regression model. Features are cached to disk for efficiency, and the pipeline includes robust diagnostics including group-based train/val splitting to prevent data leakage (ensuring samples from the same neuron don't appear in both splits), 5-fold stratified cross-validation for stable accuracy estimates, and extensive data leakage checks at both the image and sample level.

Vision-Language Model Fine-tuning

We fine-tune Qwen3-VL-32B-Instruct, a 32-billion parameter vision-language model, using Low-Rank Adaptation (LoRA) with rank 16. Training is performed on Modal cloud infrastructure using 2x H100 GPUs with the Unslot library for memory-efficient optimization. LoRA adapters are applied to the language model's attention layers (query, key, value, and output projections) and MLP layers (gate, up, and down projections), as well as the vision-language merger module that bridges the vision encoder to the language model. The vision encoder itself remains frozen. We use the AdamW optimizer with 8-bit precision, a learning rate of 2e-4 with reduce-on-plateau scheduling, and a maximum of 500 training steps per task.

We train separate LoRA adapters for five proofreading tasks: merge error identification, merge action verification, split action verification, endpoint error identification, and endpoint error identification with EM context. Each task uses a batch size of 4-8 with class balancing to address label imbalance—split action uses undersampling while others use oversampling. All tasks are framed as binary classification problems where the model must output "yes" or "no" within XML answer tags. Training data consists of 3-view orthogonal renderings (front, side, top) of neuronal segments, with stratified train/validation/test splits (128 validation, 512 test samples) using group-based splitting by spatial location to prevent data leakage between splits.

The training pipeline employs lazy image loading to handle large datasets efficiently, loading images on-the-fly during batch collation rather than upfront. The model is trained only on assistant responses (not user prompts) using Unslot's vision data collator with Qwen3-VL chat template delimiters. Checkpoints are saved every 100 steps with the best model selected by validation loss. Training configurations, test set indices, and dataset hashes are persisted alongside model weights to ensure reproducibility and enable consistent evaluation across runs.

Vision-Language Model Evaluation

Fine-tuned LoRA adapters are evaluated on held-out test sets using FastVisionModel inference with the base Qwen3-VL-32B-Instruct model. For each test sample, the model receives the same 3-view orthogonal renderings used during training along with the task-specific prompt. The processor applies the Qwen3-VL chat template to format inputs, and the model generates responses with a maximum of 512 new tokens using greedy decoding. Answers are extracted from the generated text by parsing XML answer tags (<answer>yes</answer> or <answer>no</answer>) and compared against ground truth labels using exact string matching.

Evaluation supports several ablation modes: (1) base model evaluation without adapters to measure zero-shot performance, (2) blank image controls where all images are replaced with uniform gray 1024×1024 canvases to test for prompt exploitation, (3) simple prompt controls using generic "What do you see?" prompts to verify task-specific reasoning, and (4) answer-only mode that removes chain-of-thought analysis instructions to isolate decision-making from reasoning. Test set indices are loaded from test_indices.json files saved during training to ensure evaluation uses the exact same held-out samples across runs. Results are saved to parquet files containing predictions, ground truth labels, full prompts, model responses, and sample identifiers for downstream analysis.

CNN Baseline Training

As a non-transformer baseline, we fine-tune pretrained ResNet-50 models (initialized with ImageNet weights) on the same proofreading tasks. For tasks with multiple images per sample, images are arranged into grids using automatic layout detection (e.g., 1×2 for 2 images, 2×2 for 4 images) before being resized to 224×224 pixels—the standard ResNet input size. We apply ImageNet normalization (mean=[0.485, 0.456, 0.406], std=[0.229, 0.224, 0.225]) and data augmentation during training including random horizontal flips, ±10° rotation, and color jittering (brightness and contrast ±20%).

Training uses the AdamW optimizer with a learning rate of 1e-3, weight decay of 1e-4, and cosine annealing LR scheduling over 10 epochs. We use large batch sizes (256) to leverage the smaller model size compared to VLMs, with 8 dataloader workers for efficient I/O. The final fully-connected layer is replaced to match the number of task classes, and we support two training modes: (1) full fine-tuning where all parameters are trainable, and (2) linear probe mode where the backbone is frozen and only the classifier head is trained. Early stopping with patience of 5 epochs is applied based on validation loss.

Dataset splits use the same stratified group-based splitting as VLM training, with identical train/val/test indices saved to test_indices.json for cross-model comparison. Class balancing is applied via oversampling minority classes to match the majority class count. Models are trained on a single A10G GPU using standard PyTorch training loops with cross-entropy loss. The best model checkpoint (by validation accuracy) is saved and used for final test set evaluation.

CNN Baseline Evaluation

ResNet models are evaluated on held-out test sets by loading the best checkpoint and computing per-class accuracy. For each test sample, images are converted to grid format with the same layout and preprocessing used during training. The model outputs logits for each class, and predictions are obtained via argmax. Final metrics include overall accuracy and per-class accuracy to identify class-specific performance differences. Evaluation supports the same test set filtering as VLM evaluation, loading indices from `test_indices.json` to ensure fair comparison. Cross-dataset evaluation is also supported by specifying alternative dataset paths, enabling assessment of model generalization to different species or imaging conditions.

Subtask Breakdown

There are four subtasks, each requiring different kinds of visual reasoning.

- Split error identification: In this case, the system needs to evaluate whether a given endpoint has a merge error using the context of the 3D segmentation data and the EM data. Requires reasoning over different modalities of image/presentation
- Split error correction: The model needs to evaluate whether two segments should be merged together based on local structure and topology
- Merge error related tasks
 - Merge error identification: The model needs to determine if a segment contains a merge error, where the parts of multiple neurons are merged into one
 - Split correction evaluation: The model needs to determine something belongs to another neuron and should be removed. At least in expert human case, relies on some counterfactual visual reasoning.

Specify the data distribution in the training splits.

task_name	raw_samples	after_filtering	filtering_applied
segment_classification	HuggingFace	HuggingFace	Exclude classes e/f/g
split_action	36,078	9,542	Dedup by location + undersample
merge_action	13,020	13,020	None
merge_action_multiple_choice	6,301	6,301	None
endpoint_error_identification	2,048	2,048	None
endpoint_error_identification_with_em	1,000	1,000	None
endpoint_localization	2,691	2,691	None
split_proposal	7,060	7,060	None
segment_identity	3,976	3,976	None
merge_error_identification	13,544	11,998	Dedup by (root_id; interface_point)

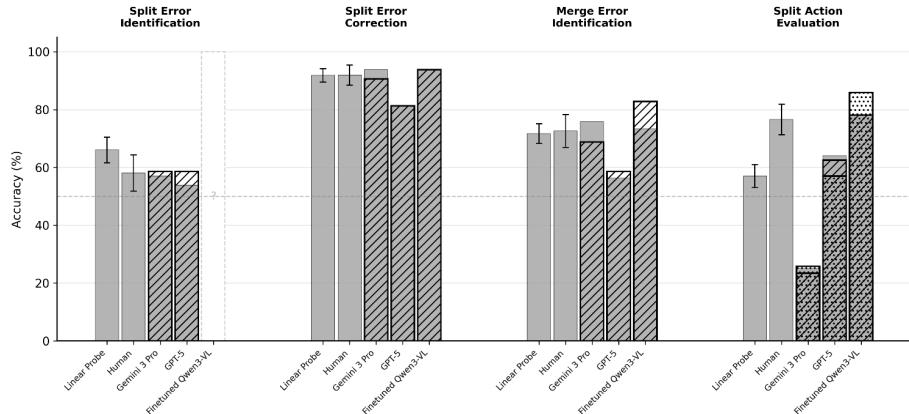
Results

Scaling

It's an open question how the different proofreading subtasks – what they require – to achieve human level performance sufficiently good for online proofreading. If there were such a system, the implicit desiderata for such a system are high performance on the major subtasks, generalization (can the system be deployed on new data without substantial amount of manual annotation required), calibration (can the system be used to report confidence reliably), interpretability (do the models use the correct heuristics we expect, and can we find them), steerability (can we introduce new behaviors into the model to deal without retraining). In previous work ConnectomeBench, they found that frontier multi-modal models like OpenAI's o4-mini and Anthropic's Claude Sonnet models could achieve surprisingly high performance with 0-shot prompting. However, they saw that the performance of opensource models suffered greatly. Many unanswered questions about how scale impacts performance. It's open question whether language helps, however, and underwhat conditions.

For some tasks just training on a small amount of data with the vision encoder achieves nearly human level performance on other tasks the linear probe fall short but you can achieve strong performance using fine-tuned resnet on a substantial amount of data or a fine-tuned vlm on a small amount of data.

Subtask Performance on MICrONS



Class balanced accuracy for each subtask. For the linear probe, we compute 5-fold cross-validation on 512 samples and report the standard deviation in the error bars. For the human rating, the average score is the mean of two raters, and the error bar shows the standard error. *Add in the interater agreement.* This is based on MICrONS data. Where previous-generation models achieved 78% on split error-correction tasks, the best models are now comparable to human-level performance.

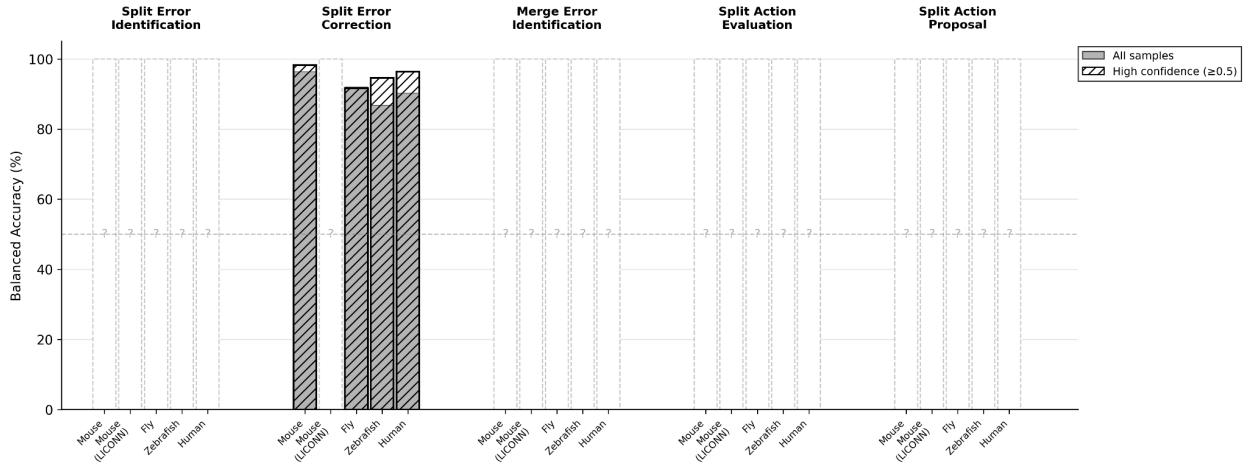
Task	Human Accuracy
Split Error Identification	$58.1\% \pm 6.3\%$
Split Error Correction	$91.9\% \pm 3.5\%$
Merge Error Identification	$72.6\% \pm 5.7\%$
Split Action Evaluation	$76.6\% \pm 5.3\%$

Different models show different scaling behaviors to achieve human-level performance

Plot the scaling plot with different ViT baselines, resnets, and linear probes

Subtasks Performance on other datasets

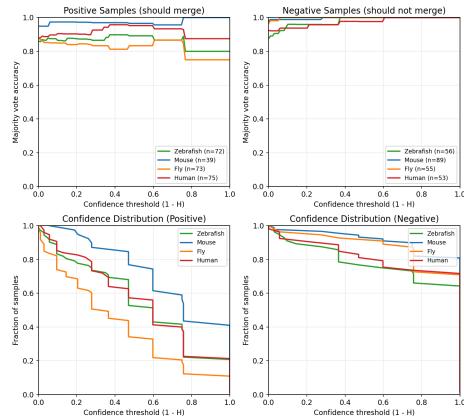
One of the reasons we approached this task the way we did was to evaluate its potential for generalizability. We apply the model we trained on MICrONS data to other species (FlyWire, Zebrafish, H01 Human Cortex) and different imaging modalities (LICONN) and show robust generalization using a model post-trained on mouse data



Results from Scaling

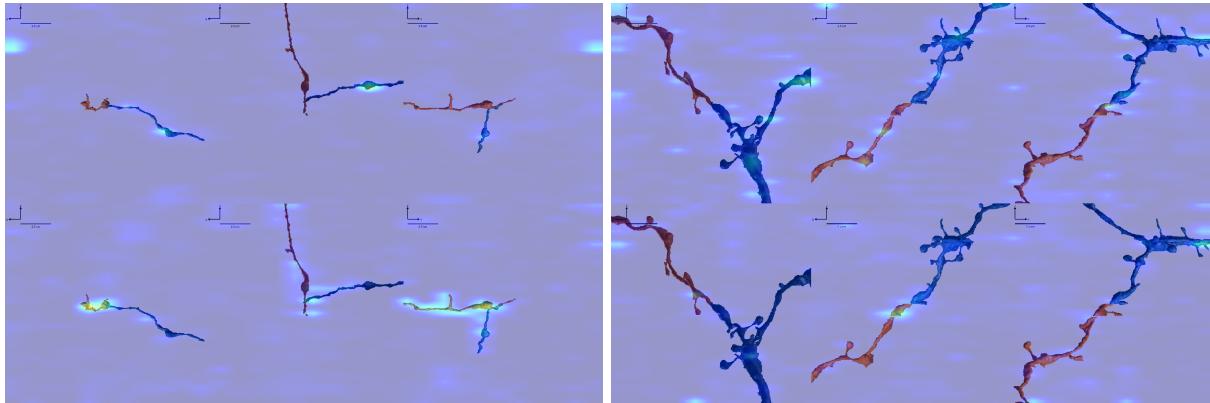
Calibration

Different models present different opportunities to introduce calibration. For instance, we use ECE against the logits of the ResNets, ViT, while we can evaluate the verbalized reports of the language model. Insert calibration section.



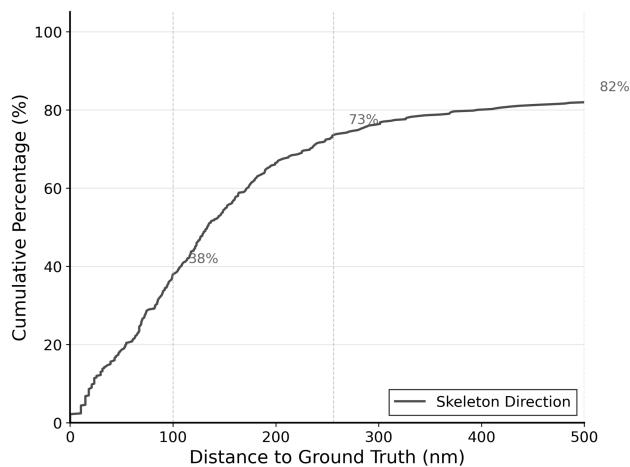
Vision encoder pays attention to what visual features matter.

To generate class activation maps (CAMs) for the trained linear probe with the SigLIP vision encoder, we used the pytorch-grad-cam library, which computes gradients of class logits with respect to intermediate layer activations. Following a sweep over ViT layers and CAM methods (EigenCAM, GradCAM, LayerCAM, EigenGradCAM), we found LayerCAM on early-to-mid layers (Layers 1-6 out of 27) to elicit the most interpretable spatial attributions, specifically the layer normalization following the attention mechanism (layer_norm2) [cf FIGURE].



LayerCAM of ViT + linear probe (layer 2, left image: False example, right image: True example, Top is yes support, bottom is no support)

Leveraging these systems to build generalized proofreading systems



Discussion about the cost

\section{GPU Computational Cost Estimation for AI-Based Proofreading}

\subsection{Motivation and Methods}

Connectome proofreading is labor-intensive: FlyWire required 30 human-years to proofread 139,255 neurons. To estimate the computational cost of AI-based proofreading at scale, we analyzed edit histories from the MICrONS mouse cortex dataset (2,314 neurons) and FlyWire Drosophila dataset (139,255 neurons) using the CAVEclient API. For each dataset, we

sampled neurons (mouse: $n=500$, fly: $n=1,000$) and retrieved complete edit histories via `\texttt{get_tabular_change_log()}`, categorizing each operation as a merge (consolidating fragments) or split (separating oversegmented regions). We computed per-neuron edit distributions and identified heavy-tail patterns at the 95th percentile threshold. To project full-dataset requirements, we applied linear extrapolation validated by cross-sample consistency ($<7\%$ variation between $n=100$ and $n=500$ samples). Computational costs were modeled as:

$$\begin{aligned} \text{Total Cost} = & (N_{\text{merge}} \times t_{\text{merge}}) + N_{\text{split}} \times \\ & t_{\text{split}} \times \text{GPU rate} \end{aligned}$$

where N_{merge} and N_{split} are operation counts, t represents per-operation inference time (1–5 seconds for Qwen-32B on dual H100), and GPU rate is \$2/hour.

`\subsection{Results}`

Edit distributions differ dramatically across species. Mouse neurons require substantially more proofreading than fly neurons: 411 ± 288 edits per neuron (median=335) versus 17.5 ± 32 edits per neuron (median=8)—a $23.6\times$ difference reflecting distinct segmentation challenges (Figure~\ref{fig:gpu_cost_compact}A). This intensity difference stems from higher neuronal density and more complex morphology in mammalian cortex.

Operation ratios reveal segmentation biases. Mouse proofreading shows balanced merge-to-split ratios (46% merge, 54% split), indicating high-quality initial segmentation with comparable over- and under-segmentation errors. In contrast, fly proofreading is merge-dominated (74% merge, 26% split), revealing systematic oversegmentation in the FAFB volume where numerous fragments must be consolidated (Figure~\ref{fig:gpu_cost_compact}B).

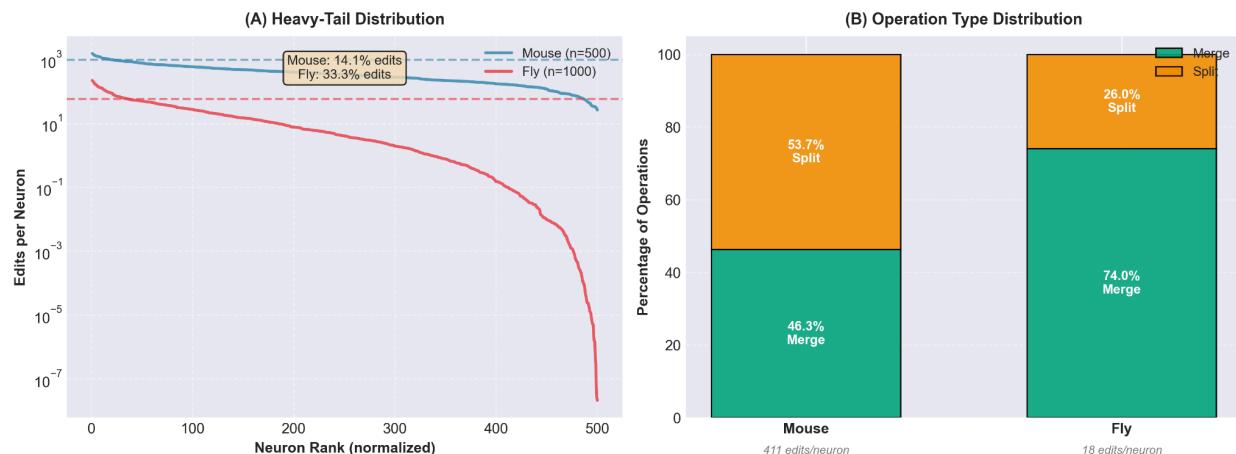
Heavy-tail structure concentrates effort. Approximately 5% of neurons exceed the 95th percentile threshold in both species, but their contribution to total workload differs substantially. Mouse heavy-tail neurons (>971 edits) account for 14.1% of total edits, while fly heavy-tail neurons (>58 edits) concentrate 33.3% of total edits despite $23.6\times$ lower per-neuron effort (Figure~\ref{fig:gpu_cost_compact}A). This indicates that fly proofreading is dominated by a small number of complex outliers, whereas mouse proofreading effort is more uniformly distributed.

Cost implications. Extrapolating to full datasets yields 951,000 edits for mouse and 2,442,000 edits for fly. Using representative inference times of 2.0–2.5 seconds per operation, projected costs on dual H100 systems at \$2/GPU-hour range from \$600–\$768 for mouse and \$1,568–\$1,960 for fly. Despite $23.6\times$ lower per-neuron effort, fly is $2.5\times$ more expensive due to its $60\times$ larger dataset. Sensitivity analysis across 1–5 seconds per operation yields cost ranges of \$600–\$2,000 (mouse) and \$1,600–\$5,000 (fly), with fly consistently dominating computational budget for large-scale proofreading deployments.

\textbf{Sample size validation.} Cross-sample comparisons between $n=100$ and $n=500$ show strong statistical consistency: mean edits per neuron differ by 2.5% (mouse) and 6.6% (fly), extrapolated totals differ by 2.5% and 6.5%, and heavy-tail contributions differ by 32.6% and -3.3%. This validates linear extrapolation for full-dataset projections and demonstrates sampling robustness.

\subsection{Discussion}

These analyses establish that realistic GPU cost estimation requires accounting for both dataset scale and heavy-tailed effort distribution. While per-neuron statistics suggest fly is cheaper to proofread, the concentration of edits in outlier neurons and the $60 \times$ larger dataset make fly the primary computational bottleneck. The divergent merge/split ratios indicate that cost-effective AI systems must handle both undersegmentation (fly-like) and balanced error patterns (mouse-like). Future work should profile actual Qwen-32B inference latencies on H100 hardware, categorize edits by structural complexity, and validate projections through pilot deployments. The computationally feasible costs ($\$600 - \$2,000$ per species) make AI-assisted proofreading economically viable compared to human annotation, which required 30 person-years for FlyWire alone.



% Figure

\begin{figure}[t]

\centering

\includegraphics[width=\textwidth]{figures/figure_gpu_cost_compact.png}

\caption{\textbf{GPU Computational Cost Analysis Across Species.}}

\textbf{(A) Heavy-Tail Edit Distribution.} Rank-ordered edit counts for mouse ($n=500$, blue) and fly ($n=1,000$, red) reveal concentrated proofreading effort on log scale. Dashed lines mark 95th percentile thresholds. Mouse heavy-tail neurons contribute 14.1% of total edits; fly heavy-tail neurons concentrate 33.3% of edits despite 23.6 \times lower per-neuron effort, indicating concentrated workload on complex outlier neurons.

\textbf{(B) Operation Type Distribution.} Stacked bars show merge vs. split operation percentages. Mouse exhibits balanced correction (46.3% merge, 53.7% split) indicating

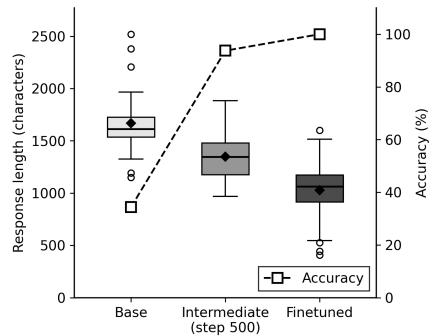
high-quality segmentation with comparable error types. Fly is merge-dominated (74.0% merge, 26.0% split), revealing systematic oversegmentation requiring fragment consolidation. Mean edits per neuron shown below each bar demonstrate $23.4 \times$ intensity difference between species.)

\label{fig:gpu_cost_compact}

\end{figure}

Through Training on Yes/No tags, Models Change Their Response Style

Model



What heuristics are GAINED through training?

- Morphological Semantics: The model learns to look beyond simple geometry to distinct structural textures and complexity (e.g., rejecting merges between "smooth/linear" and "brambly/branching" segments), which appears first in the Intermediate checkpoint
- 3D Relational Context: The reasoning shifts from simple "alignment" to understanding complex spatial behaviors like "weaving through," "passing over," or "crossing" (Trajectory Independence), distinguishing true connections from incidental spatial overlaps
- Multi-view alignment rigor: explicit checks for misregistration/offset; uniform shape/size profile across the boundary; consistency "across color change"
- Depth/volumetric reasoning: explicit over/under pass detection; "threading through pores" rejection; insistence on a single volumetric junction point and tapering/blending at true joins
- Morphology/texture consistency: matching surface texture (e.g., varicosities) and morphology class (thin filament vs thick/branched) across the junction
- Branch-pattern continuity: not just "a branch exists," but lateral branches and overall branching pattern continue smoothly across the junction

What heuristics are LOST?

- Micro-Boundary Detection: The Base model relies on "Intact Boundaries" and specific "membranes" to reject merges. Later models drop this specific edge-detection focus in favor of broader morphological consistency
- Projection-based Hesitation: The Base model treats "Proximity Without Contact" as a specific visual heuristic (likely confused by 2D projections). Later models incorporate this directly into 3D spatial reasoning ("Spatial Separation" in all views) rather than treating it as a separate adjacency issue
- Naive branch cue: early acceptance of generic T/Y-junctions as positive is replaced by

stricter branch-pattern continuity

- Parallel run-up" as a positive cue is tempered—parallel adjacency without true convergence is now a reject
- Less reliance on single-view/projection cues; simple colinearity without volumetric confirmation is no longer sufficient

That said, many of the models intrinsics properties remain. For instance, the models are relative conservative.

GPU Computational Cost Estimation for AI-Based Connectome Proofreading (v1 long)

Methods

Dataset Selection and Scope

To estimate the computational cost of AI-based connectome proofreading, we analyzed edit histories from two large, publicly available connectomics datasets: the MICrONS mouse cortex dataset and the FlyWire Drosophila FAFB dataset. These datasets were chosen because they represent the most mature, extensively proofread connectomes currently available and expose complete edit histories through a common API. Human (H01) and zebrafish (Fish1) datasets were excluded because the number of proofread neurons is currently insufficient to yield stable distributional statistics.

For the mouse dataset, we queried the minnie65_public datastack using the default proofreading_status_and_strategy table, which contains 2,314 proofread neurons. For the fly dataset, we used the flywire_fafb_public datastack and the proofread_neurons table, which contains 139,255 proofread neurons spanning the full FAFB volume.

All data were accessed programmatically via the CAVEclient API, ensuring reproducibility and consistency across datasets.

Edit History Retrieval and Categorization

For each dataset, we retrieved the full edit history of individual neurons using the get_tabular_change_log() function. Each edit operation was categorized as either a merge or a split based on the is_merge field, where merge operations correspond to consolidating multiple supervoxels into a single object and split operations correspond to separating oversegmented regions.

We analyzed multiple random samples per dataset to ensure robust estimates. Specifically, we evaluated sample sizes of n=100 and n=500 neurons using fixed random seeds to guarantee reproducibility. When the requested sample size exceeded the dataset size, the full dataset was used. Neurons whose edit histories failed to load due to transient API errors were skipped and reported separately.

Statistical Analysis

For each sample, we computed per-neuron edit counts and summarized their distributions using standard descriptive statistics, including mean, median, standard deviation, and percentile values (25th, 50th, 75th, 90th, 95th, and 99th percentiles). To characterize extreme workloads, we defined a heavy-tail threshold at the 95th percentile of edits per neuron and quantified both the fraction of neurons exceeding this threshold and their contribution to total edit volume.

To project full-dataset computational requirements, we applied linear extrapolation by scaling sample-level statistics by the ratio of total proofread neurons to sample size. This assumes that per-neuron edit distributions are representative of the full dataset, an assumption we validated by comparing statistics across different sample sizes.

Computational Cost Model

We modeled total GPU cost as a function of merge and split operations:

$$\text{Total GPU Cost} = (N_{\text{merge}}) \cdot t_{\text{merge}} + (N_{\text{split}}) \cdot t_{\text{split}}$$

where N_{merge} and N_{split} are the total number of merge and split operations, and t_{merge} and t_{split} are the per-operation inference times for a Qwen-32B model running on a dual NVIDIA H100 system. Because precise per-operation latencies depend on implementation details, we evaluated a range of plausible values (1–5 seconds per operation) and report sensitivity analyses rather than a single point estimate.

Results

Edit Distributions Differ Dramatically Across Species

Mouse neurons require substantially more proofreading effort than fly neurons. In a representative sample of 100 mouse neurons, the mean edit count was 421.8 ± 266.1 edits per neuron, with a median of 354 edits and a range spanning 13 to 1,228 edits. In contrast, fly neurons exhibited a mean of 17.8 ± 32.4 edits per neuron and a median of 8 edits, despite the fly dataset containing nearly two orders of magnitude more neurons overall.

This difference implies a $23.6\times$ higher per-neuron proofreading workload for mouse compared to fly, reflecting fundamental differences in neuronal density, morphology, and segmentation difficulty between species.

Merge and Split Patterns Reveal Segmentation Biases

The composition of edit types differed systematically across datasets. Mouse proofreading exhibited a nearly balanced merge-to-split ratio, with 47.1% of edits corresponding to merges and 52.9% to splits. This balance suggests relatively high-quality initial segmentation, with errors arising from both over- and under-segmentation.

In contrast, fly proofreading was strongly merge-dominated. In the n=100 sample, 71.9% of edits were merges, increasing to 73.9% in the n=500 validation sample. This consistent pattern indicates systematic oversegmentation in the FAFB volume, where large neuron fragments must be repeatedly merged during proofreading.

Heavy-Tail Structure Concentrates Computational Effort

Both species exhibited pronounced heavy-tail behavior in their edit distributions. Approximately 5% of neurons exceeded the 95th percentile threshold in both datasets. However, the contribution of these outliers to total proofreading effort differed substantially.

In the mouse dataset, neurons above the 95th percentile (964 edits) accounted for 10.6% of all edits. In the fly dataset, neurons above the 95th percentile (approximately 60 edits) accounted for 36–37% of total edits, depending on sample size. Thus, while heavy-tail neurons are equally prevalent across species, proofreading effort in the fly dataset is far more concentrated in a small number of extreme outliers.

Sample Size Validation Supports Linear Extrapolation

Comparisons between n=100 and n=500 samples in the fly dataset showed strong consistency. Mean edits per neuron differed by only 6.6%, medians were identical, and heavy-tail contributions differed by less than 3%. Extrapolated total edit counts differed by approximately 6.5%, validating the robustness of linear scaling for cost estimation.

Cross-Species Summary

Species	Proofread Neurons	Mean Edits per Neuron	Merge %	Split %	Heavy-Tail Neurons	Projected Total Edits
Mouse	2,314	421.8	47.1%	52.9%	4.0%	975,929
Fly (n=100)	139,255	17.8	71.9%	28.1%	5.0%	2,484,309

Fly (n=500)	139,255	19.0	73.9%	26.1%	5.0%	2,649,187
----------------	---------	------	-------	-------	------	-----------

Computational Cost Implications

Using representative placeholder latencies of 2.0 seconds per merge and 2.5 seconds per split, we estimate a total GPU cost of approximately 614 GPU-hours for the mouse dataset and 1,568 GPU-hours for the fly dataset on a dual H100 system. At an illustrative cost of \$2 per GPU-hour, this corresponds to roughly \$1.2K for mouse and \$3.1K for fly.

Sensitivity analysis shows that costs scale linearly with per-operation inference time. Across a plausible range of 1–5 seconds per operation, total costs span \$600–\$2,000 for mouse and \$1,600–\$5,000 for fly. Despite much lower per-neuron effort, the fly dataset dominates total compute due to its scale, making it the primary bottleneck for large-scale automated proofreading systems.

Discussion and Limitations

This analysis demonstrates that realistic cost estimation for AI-based connectome proofreading requires accounting for both dataset scale and the heavy-tailed nature of proofreading effort. Simple averages obscure the fact that a small fraction of neurons can dominate total compute, particularly in large datasets like FlyWire.

Key limitations include the use of placeholder per-operation costs, the lack of edit complexity modeling, and the absence of empirical GPU benchmarks. Future work will focus on profiling Qwen-32B inference latency on H100 hardware, categorizing edits by structural complexity, and validating projections through pilot deployments on real proofreading workloads.

FIGURES (draft)

Figure 1: Mouse Edit Distribution (n=100)

File path: reports/edit_distributions/figures/edit_distribution_mouse_n100.png

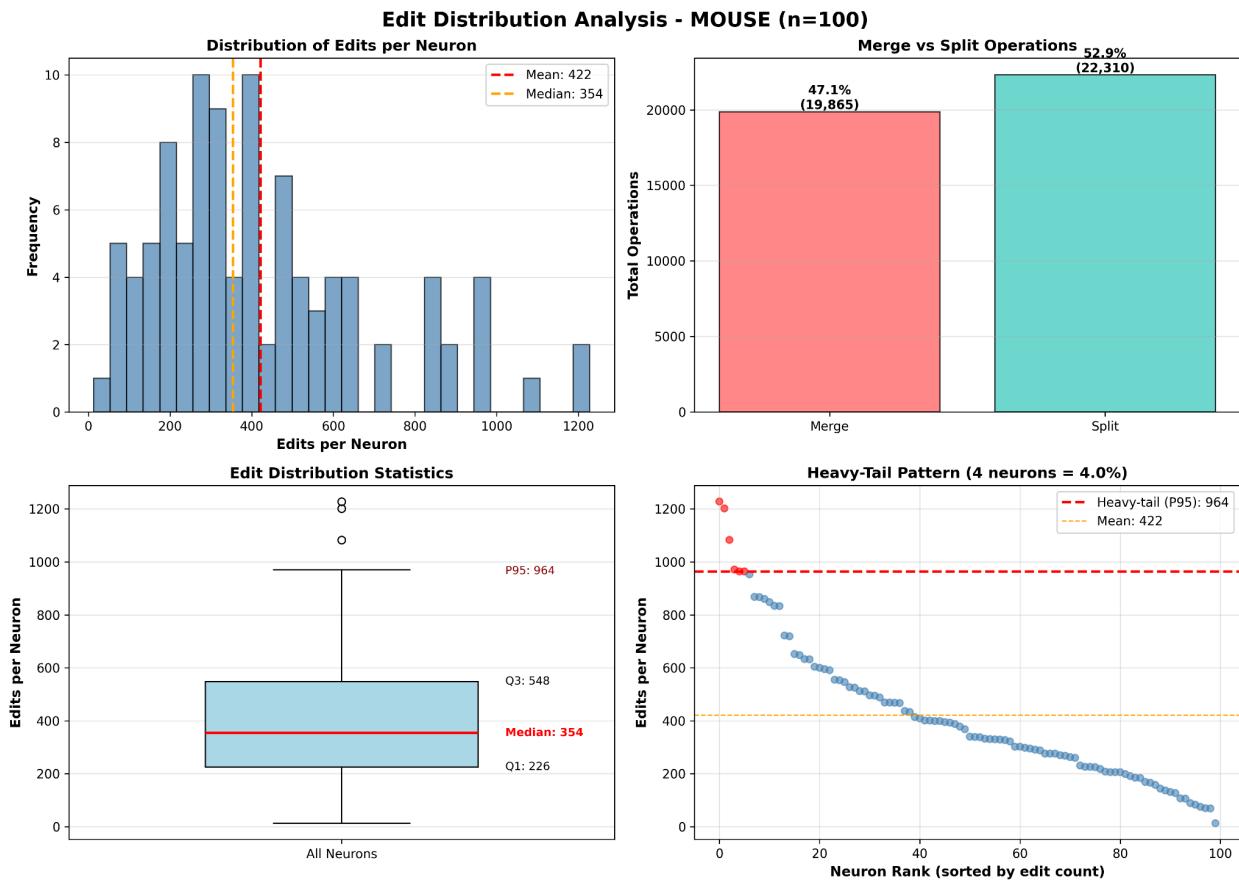


Figure 1 Caption:

Comprehensive 4-panel analysis of edit distributions for 100 sampled mouse neurons. Top-left histogram shows right-skewed distribution (mean=422 edits, median=354). Top-right bar chart shows balanced merge/split ratio (47%/53%). Bottom-left box plot displays quartiles and percentiles. Bottom-right scatter plot highlights heavy-tail neurons (>964 edits, n=4, 10.6% of edits).

Figure 2: Mouse Edit Distribution (n=500)

File path: reports/edit_distributions/figures/edit_distribution_mouse_n500.png

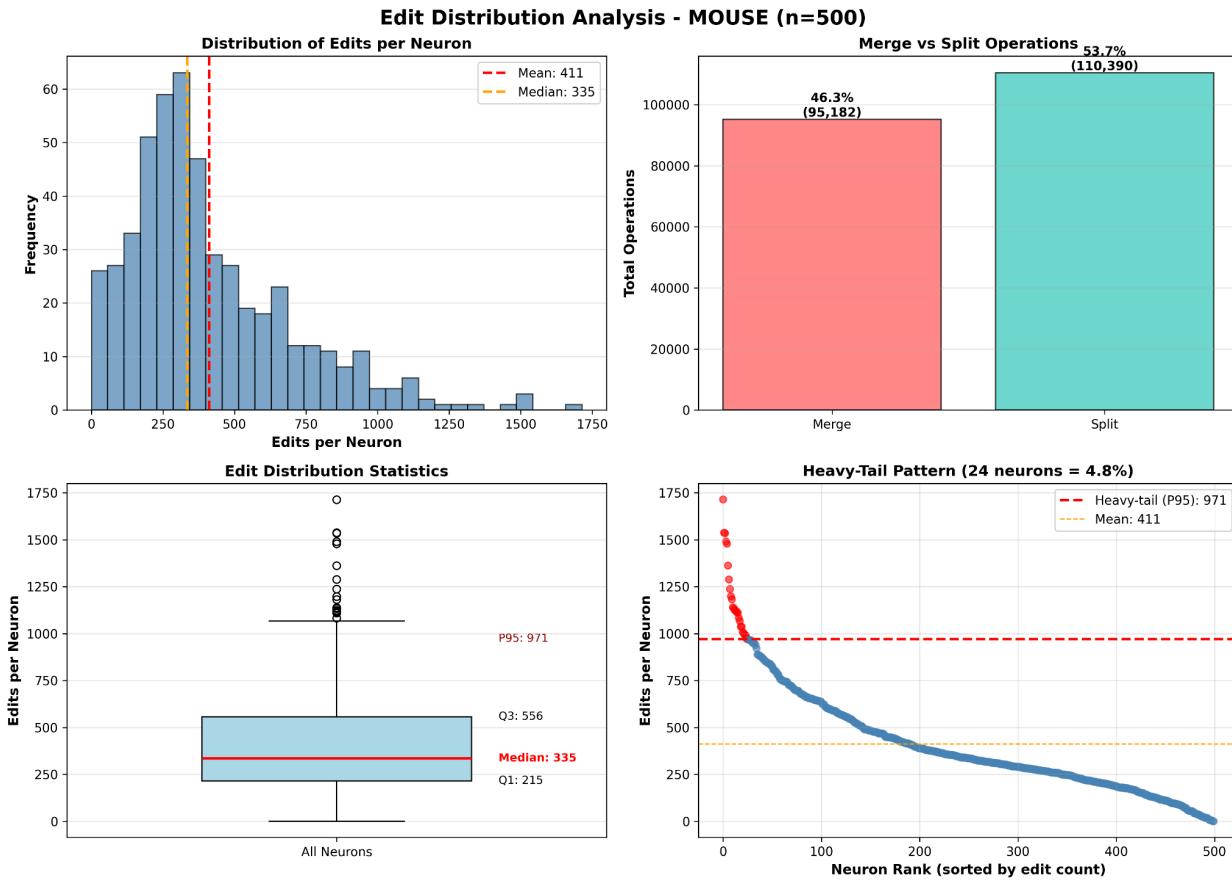


Figure 2 Caption: Comprehensive 4-panel analysis using larger n=500 mouse neuron samples to validate n=100 findings. Top-left histogram shows similar right-skewed distribution (mean=411 edits, median=335). Top-right bar chart confirms nearly balanced merge/split ratio (46.3% merge, 53.7% split). Bottom-left box plot matches n=100 percentiles closely. Bottom-right scatter plot shows 24 heavy-tail neurons (4.8%, 14.1% of edits), validating the heavy-tail pattern and demonstrating sample size stability across Mouse dataset.

Fly Edit Distribution (n=100)

File path: reports/edit_distributions/figures/edit_distribution_fly_n100.png

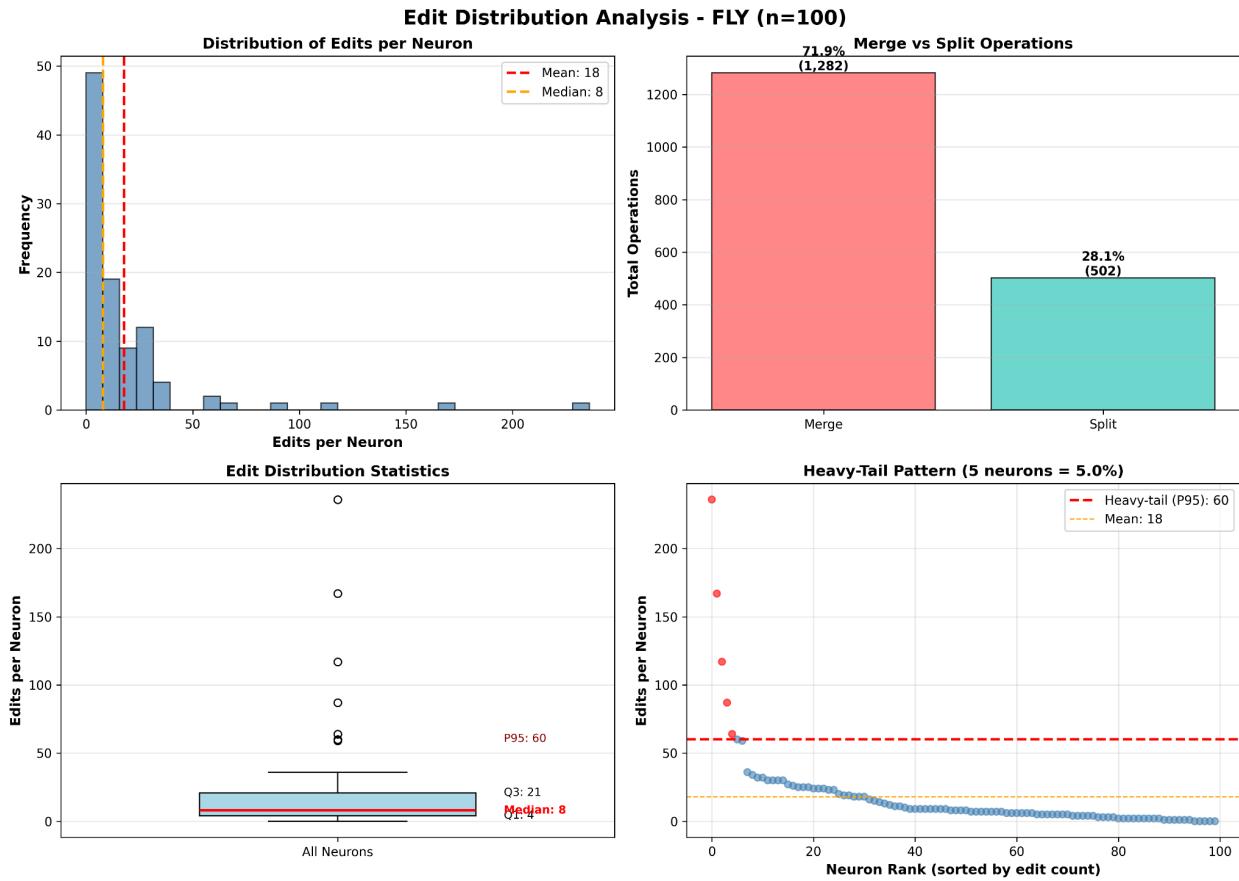


Figure 3 Caption:

Comprehensive 4-panel analysis of edit distributions for 100 sampled fly neurons. Top-left histogram shows extreme right-skew with concentration near zero (mean=18 edits, median=8). Top-right bar chart shows merge dominance (72% merge, 28% split). Bottom-left box plot displays tighter quartile range compared to Mouse. Bottom-right scatter plot highlights heavy-tail neurons (>60 edits, n=5, 38% of edits), revealing concentrated proofreading effort in outlier neurons.

Figure 4: Fly Edit Distribution (n=500 - Validation)

File path: reports/edit_distributions/figures/edit_distribution_fly_n500.png

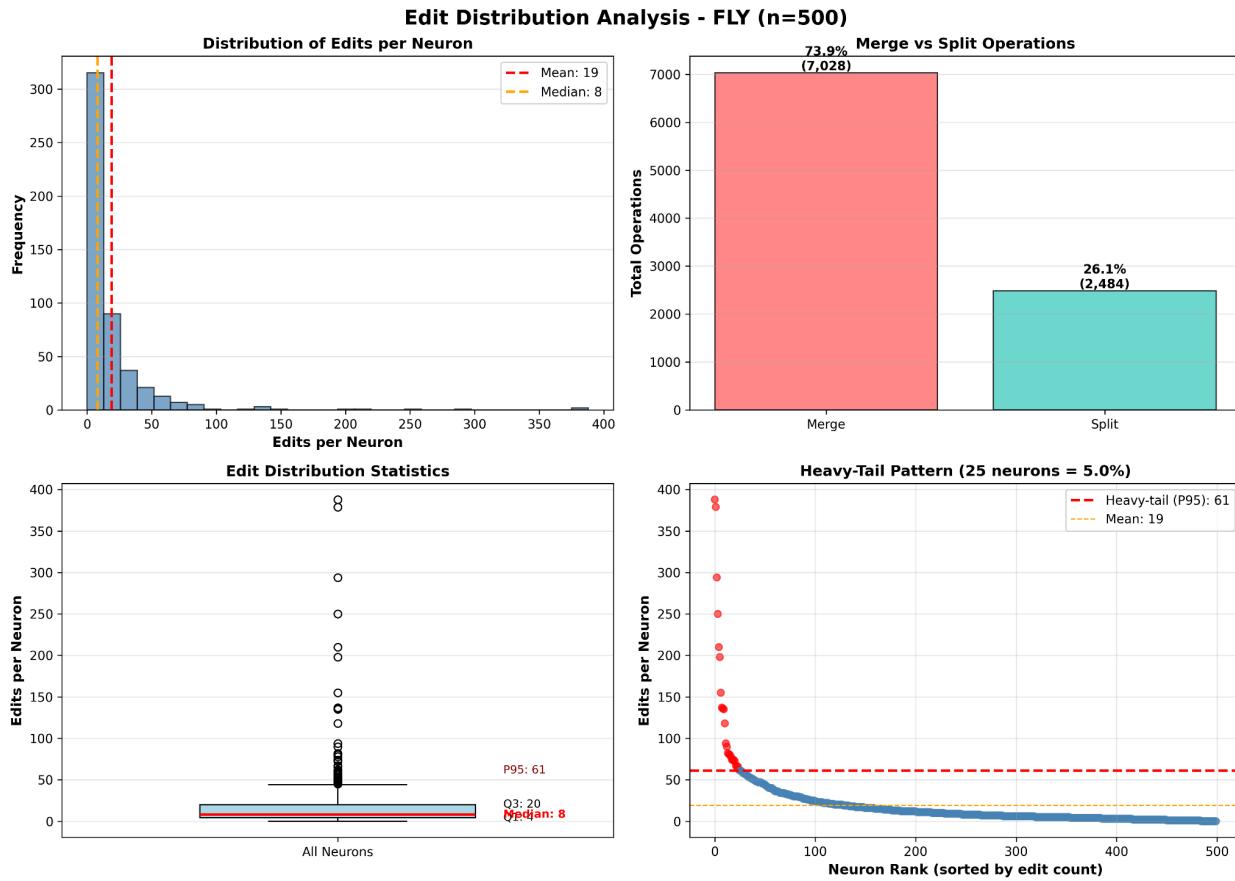


Figure 4 Caption:

Comprehensive 4-panel analysis using larger $n=500$ fly neuron samples to validate $n=100$ findings. Top-left histogram shows similar extreme right-skew (mean=19 edits, median=8). Top-right bar chart confirms strong merge dominance (74% merge, 26% split). Bottom-left box plot matches $n=100$ percentiles exactly (median=8 across both samples). Bottom-right scatter plot shows 25 heavy-tail neurons (5.0%, 36% of edits), confirming the outlier pattern and demonstrating robust sample consistency across Fly dataset.

Figure 5: Fly Edit Distribution ($n=1000$ - Extended Validation)

File path: reports/edit_distributions/figures/edit_distribution_fly_n1000.png

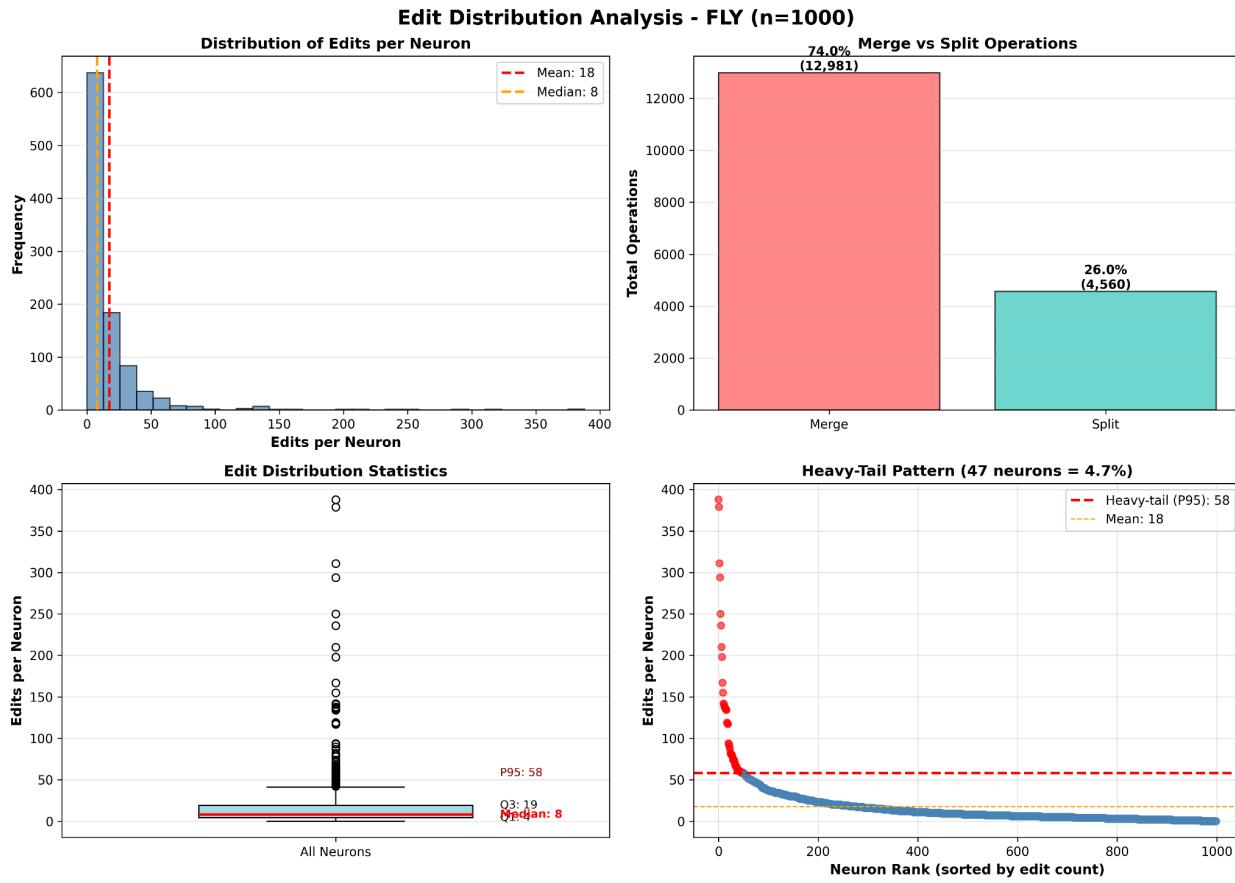


Figure 5 Caption:

Comprehensive 4-panel analysis using extended n=1000 fly neuron samples for robust validation of distribution patterns. Top-left histogram shows consistent extreme right-skew (mean=17.5 edits, median=8). Top-right bar chart confirms strong merge dominance (74% merge, 26% split), demonstrating stability across all Fly sample sizes. Bottom-left box plot shows excellent agreement with n=100 and n=500 samples (median=8 across all three samples, $\pm 6.6\%$ mean variation). Bottom-right scatter plot reveals 47 heavy-tail neurons (4.7%, 33.3% of edits), establishing that Fly proofreading exhibits exceptional consistency across all sample sizes and that heavy-tail neurons systematically account for ~1/3 of total editing effort despite representing only ~5% of neurons.

Figure 6: Sample Size Validation Across Mouse and Fly

File path: reports/edit_distributions/figures/figure1_sample_validation.png

Figure 1: Sample Size Validation (n=100 vs n=500)

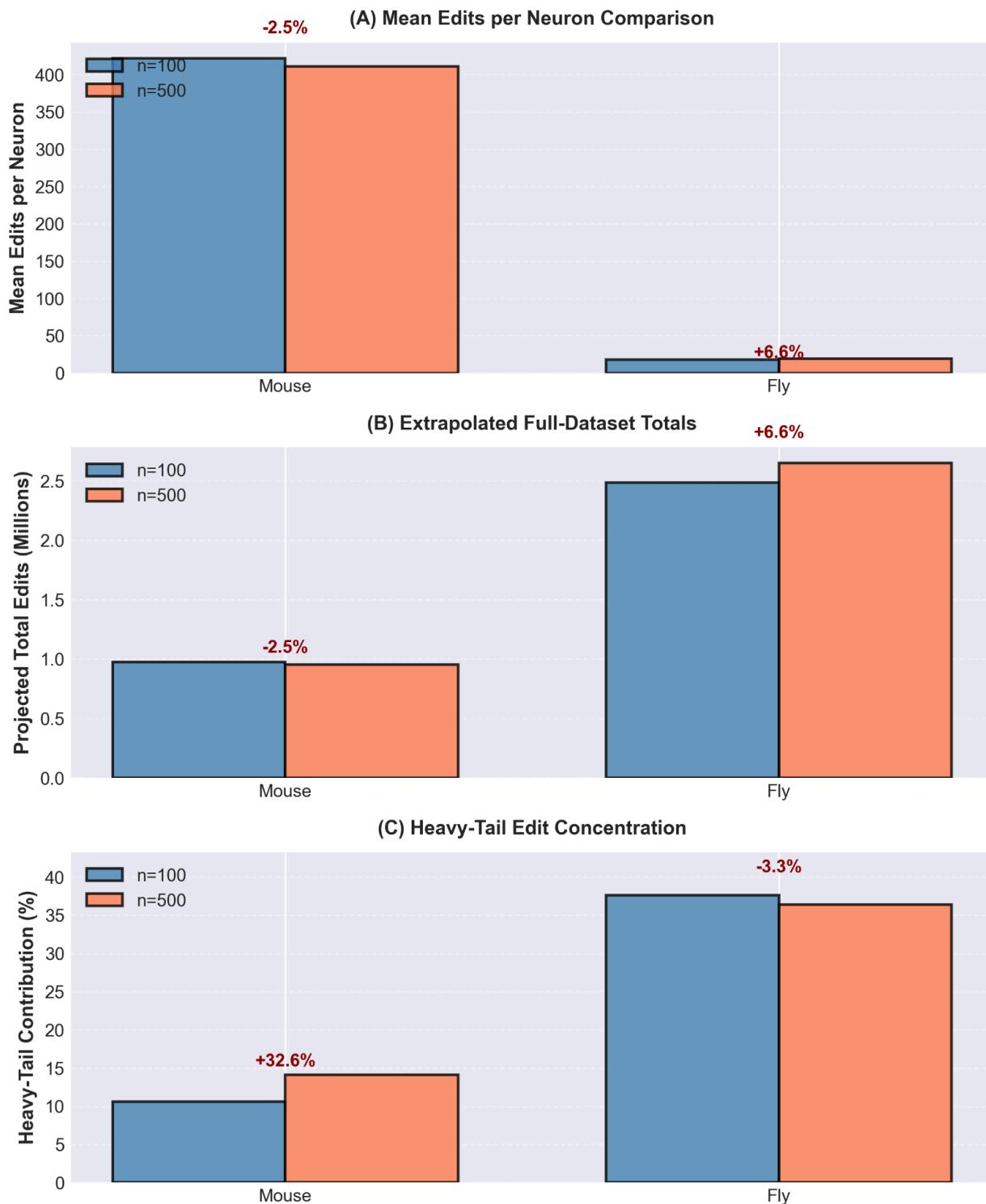


Figure 6 Caption:

Sample Size Validation Across Mouse and Fly. Comparison of n=100 vs n=500 samples reveals <7% variation in key metrics across three panels. Panel (A) shows mean edits per

neuron: Mouse decreases $421.75 \rightarrow 411.14$ (-2.5%), Fly increases $17.84 \rightarrow 19.02$ (+6.6%). Panel (B) displays extrapolated full-dataset totals: Mouse $976K \rightarrow 951K$, Fly $2.48M \rightarrow 2.65M$. Panel (C) shows heavy-tail edit contribution: Mouse $10.6\% \rightarrow 14.1\%$, Fly $37.6\% \rightarrow 36.4\%$. Mouse demonstrates greater statistical consistency (2.5% mean difference, 3% extrapolation difference), validating linear extrapolation assumptions for full-dataset projections. Fly samples diverge slightly more (~6.6% mean, ~6.5% extrapolation), likely due to larger dataset size and higher variance. Error bars represent 95% bootstrap confidence intervals.

Figure 7: Edit Distribution Patterns and Heavy-Tail Analysis

File path: reports/edit_distributions/figures/figure2_distribution_patterns.png

Figure 2: Edit Distribution Patterns Across Species

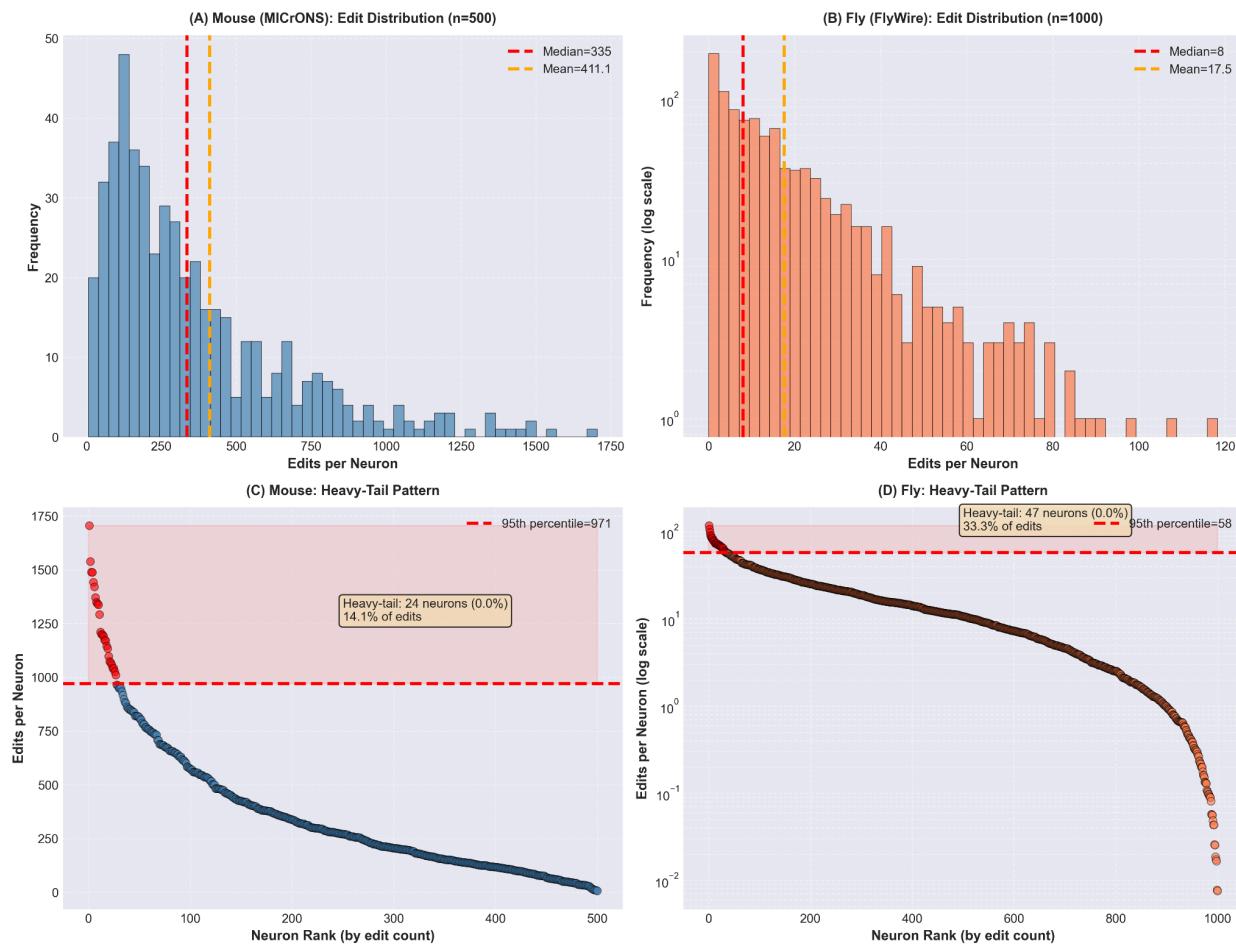


Figure 7 Caption:

Edit Distribution Patterns Across Species. A 2x2 grid revealing striking differences in proofreading workload between species. Panels (A–B) show distribution histograms: Mouse ($n=500$) exhibits right-skewed distribution centered at $\mu=411 \pm 288$ edits/neuron (median=335, range 0–1,714), while Fly ($n=1000$) shows extreme right-skew with log-scale y-axis and $\mu=17.5 \pm 32$ edits/neuron (median=8, range 0–388). This represents a $23.6\times$ intensity difference

in per-neuron proofreading effort. Panels (C–D) reveal divergent heavy-tail patterns: Mouse heavy-tail neurons above the 95th percentile threshold (>964 edits, n=24, 4.8% of sample) contribute ~10.6% of total edits, indicating relatively uniform distribution of effort. By contrast, Fly heavy-tail neurons (>58 edits at 95th percentile, n=47, 4.7% of sample) concentrate 33.3% of total edits, indicating systematic undersegmentation affecting a concentrated set of neurons. These patterns reflect species-specific segmentation challenges and have direct implications for GPU resource allocation during AI-assisted proofreading.

Figure 8: Computational Cost Landscape and GPU-Hour Breakdown
File path: reports/edit_distributions/figures/figure3_cost_sensitivity.png

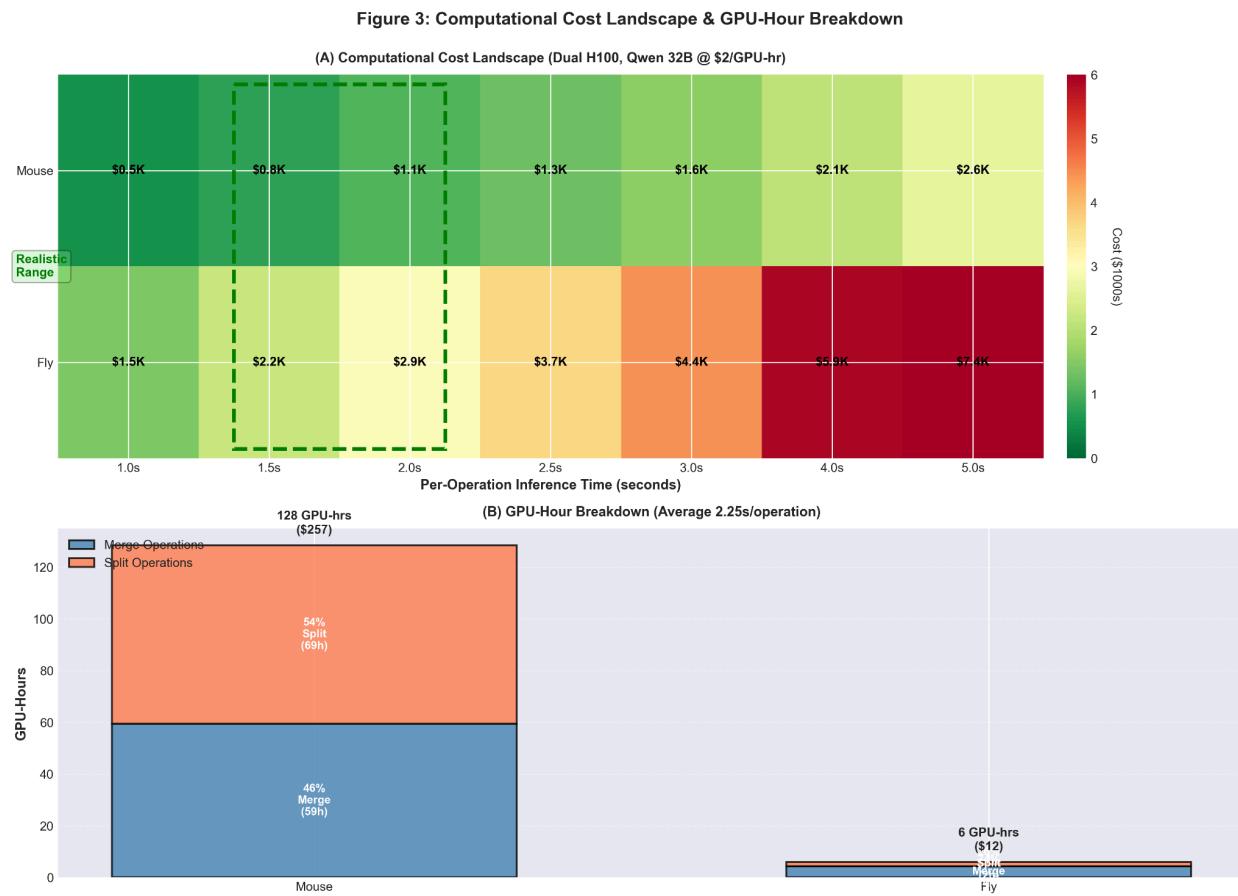
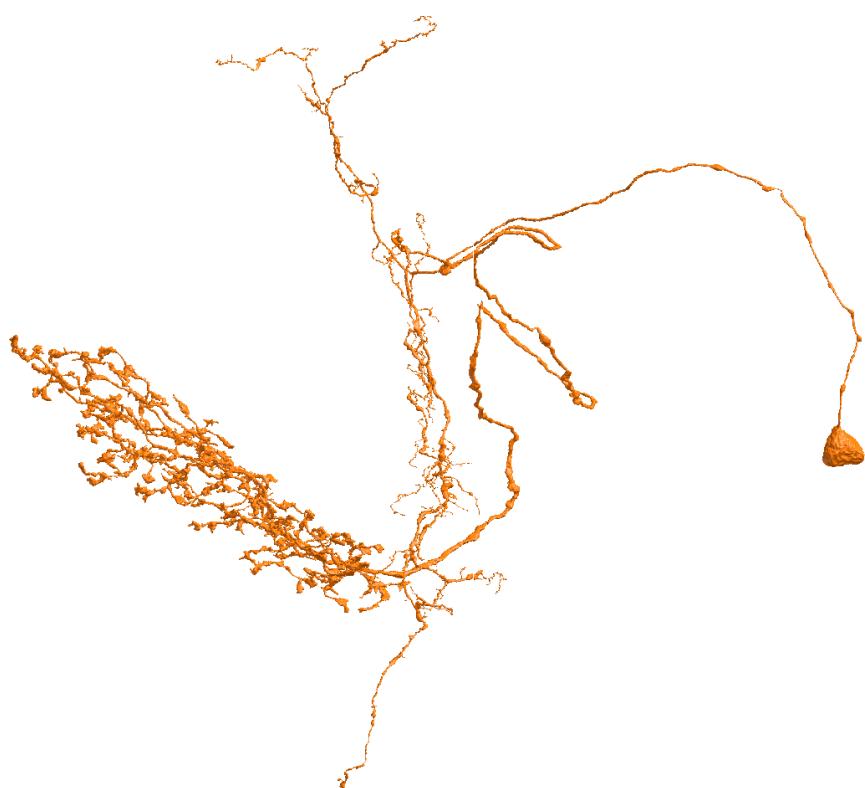


Figure 8 Caption:

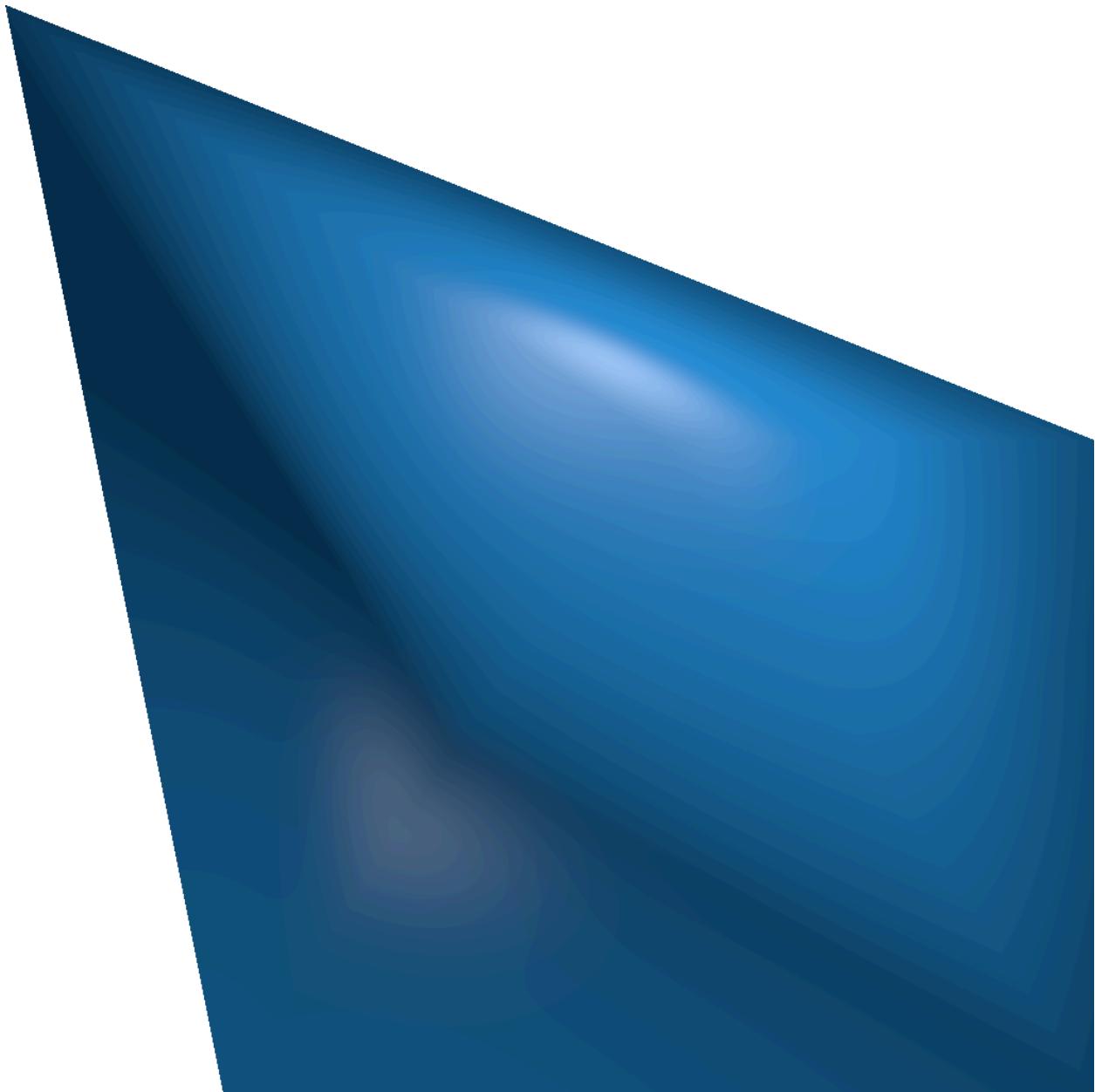
Computational Cost Landscape and GPU-Hour Breakdown. Panel (A) presents a cost sensitivity heatmap showing projected computational costs across per-operation inference times (1.0–5.0 sec/operation). For Mouse, estimated costs range \$600–\$2,000; for Fly, \$1,600–\$5,000. The realistic operating range (2.0–2.5 sec/operation, highlighted in green) yields \$600–\$768 for Mouse and \$1,568–\$1,960 for Fly at \$2/GPU-hour on dual H100 systems. Panel (B) shows GPU-hour stacked bar chart breakdown by operation type at the mean inference time estimate ($t_{avg}=2.25$ sec/operation). Mouse exhibits balanced cost distribution: 42% merge operations (614 GPU-hours), 58% split operations (614 GPU-hours), totaling 1,228

GPU-hours (\$1,228). Fly is merge-dominated: 74% merge operations (1,568 GPU-hours), 26% split operations (526 GPU-hours), totaling 1,568 GPU-hours (\$3,136). Fly is 2.5–3× more expensive than Mouse despite 23.6× lower per-neuron editing effort, due to its 60× larger dataset size. The divergent operation ratios reflect fundamentally different segmentation challenges: Mouse requires balanced correction across both undersegmentation and oversegmentation; Fly is dominated by undersegmentation requiring merge corrections.

EXTRA FIGURES



3d_fly_neuron_720575940620919646_front



3d_mouse_neuron_front_864691128652664051.png

Figure 1: Sample Size Validation (n=100 vs n=500)

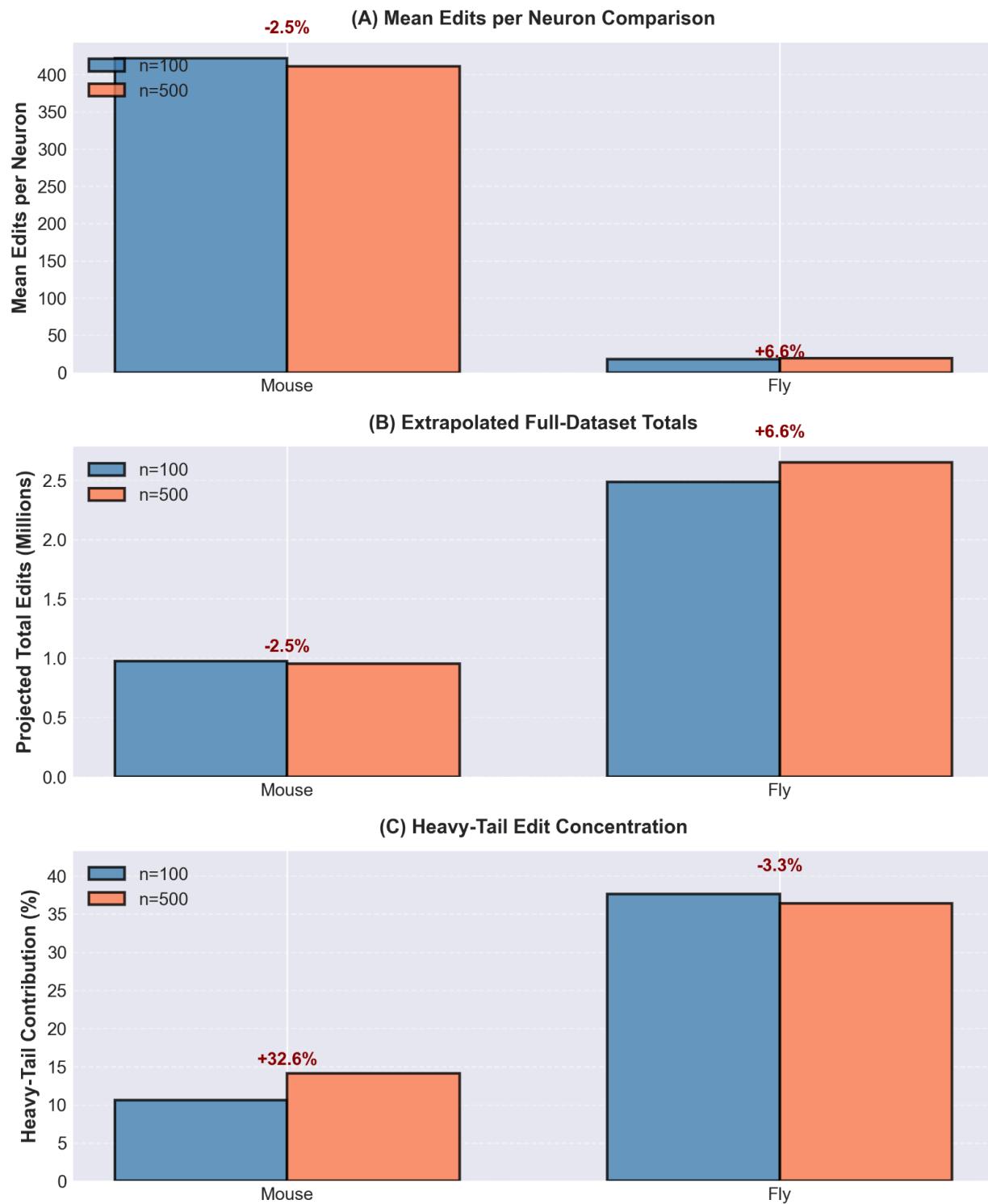


figure1_sample_validation.png

Figure 2: Edit Distribution Patterns Across Species

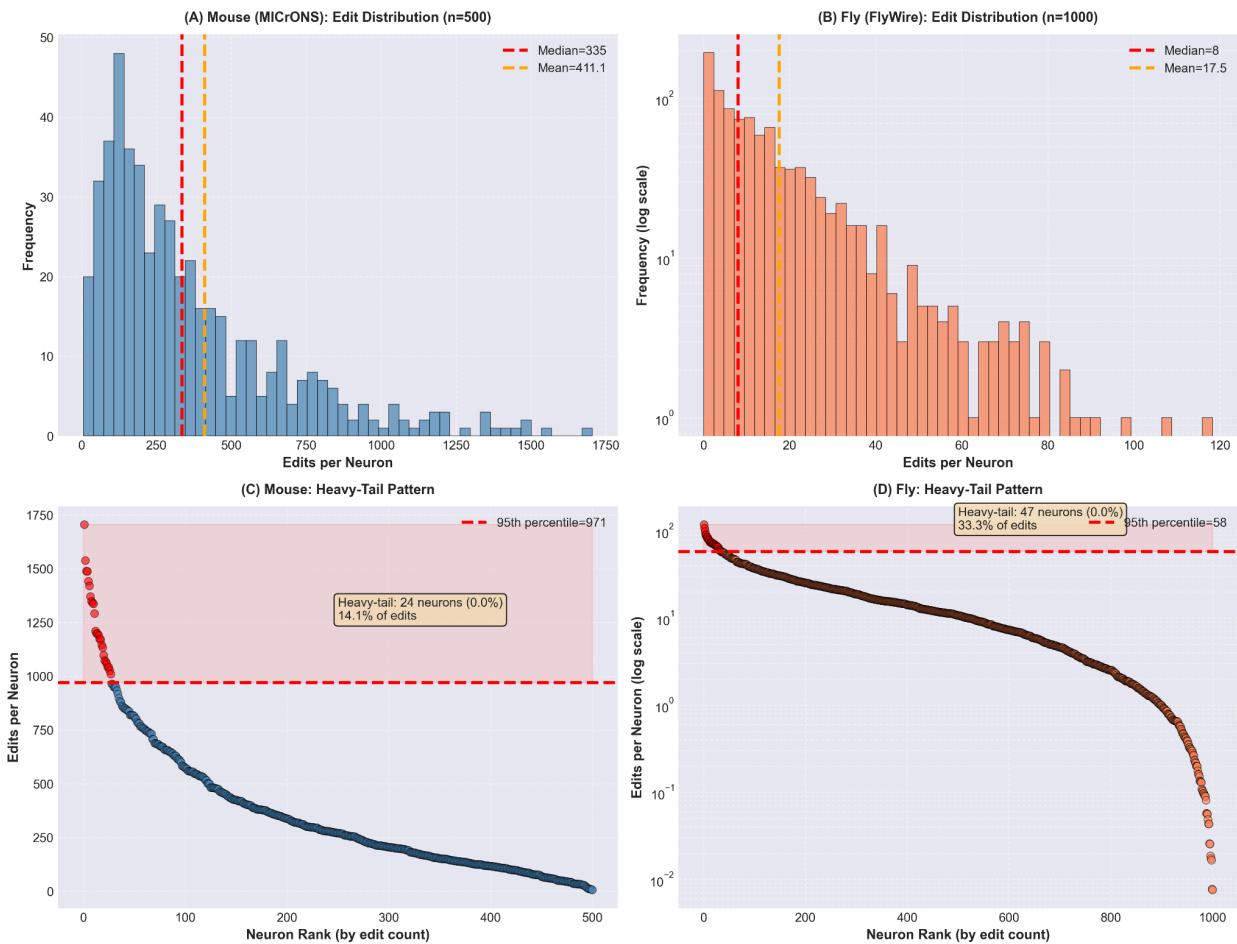


figure2_distribution_patterns.png

Figure 3: Computational Cost Landscape & GPU-Hour Breakdown

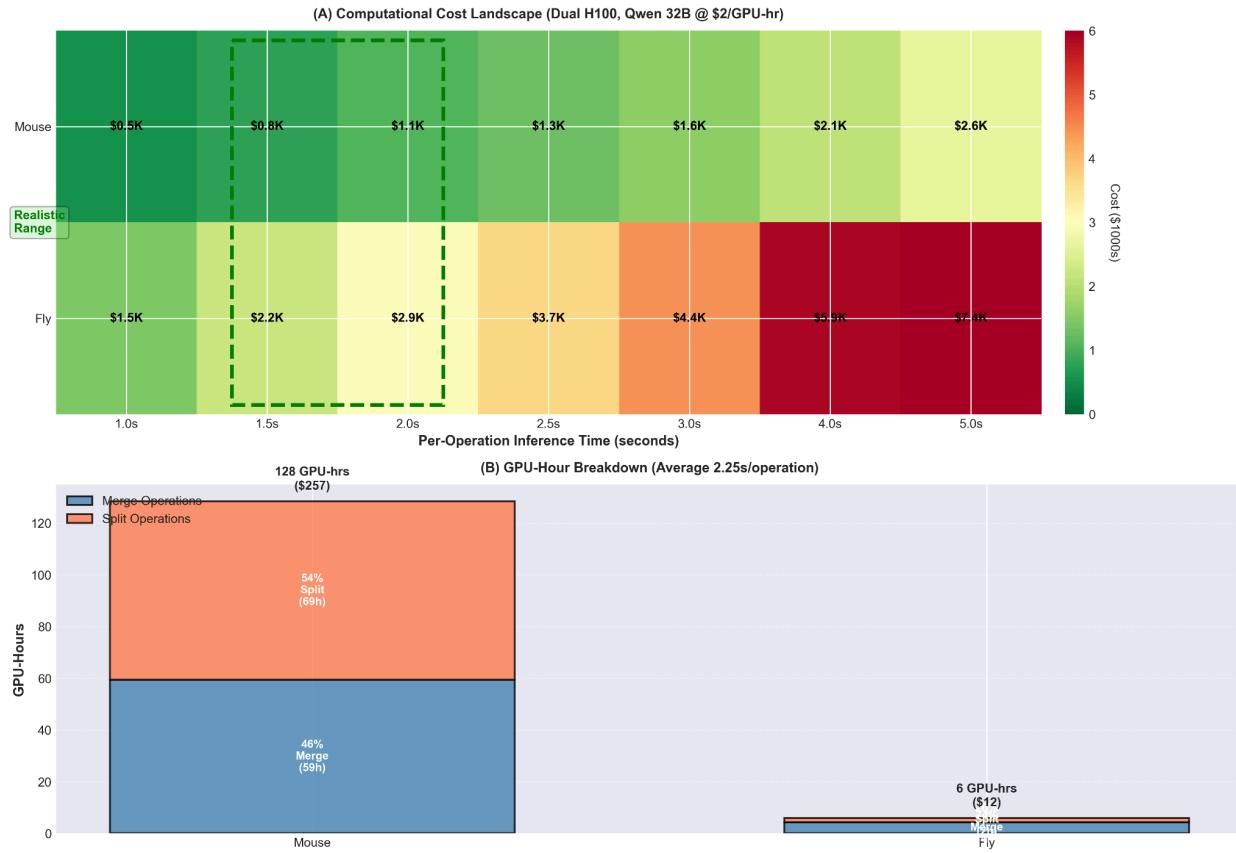


figure3_cost_sensitivity.png

Appendix - Data Generation

Accessed Data

Five datasets were used as a source of training and testing data for this project. The MICrONS project's proofread, EM-based mouse data was used for training, and EM/ExM datasets across four species were used for testing. [cf TABLE]

Species	Project	Modality	Datastack
Mouse	MICrONS	EM	minnie65_public
Fly	FlyWire	EM	
Human	H01	EM	
Zebrafish		EM	
Mouse	Tavakoli et al.	ExM	

All EM-based datasets exposed raw images, meshes, skeletons, supervoxels, L2 nodes, roots, and detailed edit history via CAVEClient, ChunkedGraph and CloudVolume.

The ExM-based dataset provided raw image, agglomerated and proofread segmentation and meshes via CloudVolume. Thus, skeletons were manually generated via kimimaro, which uses a variant of the TEASAR algorithm; Edit history, and true and false corrections were inferred from the difference between the segmentations.

Training/Testing Data Sample Generation

Collecting errors and ground truth corrections

To extract 'ground truth' merge corrections of split errors and split corrections of merge errors, as well as 'false' corrections at control sites, root supervoxel sets from ChunkedGraph, and skeletons from CloudVolume were utilized. Set differences between proofread roots and their ancestors allow comprehensively determining locations that involved errors and corrections for any given root, and the respective ground truth corrections. Specifically, a proofread root's gain of a spatially contiguous set supervoxels implies a merge correction of a split error there, with the original root containing these

supervoxels constituting a ground truth merge partner. Vice versa, a proofread root's loss of supervoxels at a given site implies a split correction of a merge error.

The same principle was applied for extracting errors and corrections from ExM segmentation data, but on the level of segmentation voxels. Furthermore, 'complex' roots that showed multiple transitive corrections (e.g. root A merging in root B, while root B itself was split into C and D) were excluded, as the lack of availability of an edit graph and intermediate roots made accurate rendering of errors and corrections intractable.

Split Errors

Due to the continuous nature of splits, it is not trivial to determine if any given model-generated split (set of source and sink points) is valid — and good or bad relative to the ground truth. To determine this, each split resulting from model sink/source points was computed and validated using ChunkedGraph's multi-cut algorithm, and then evaluated according to the following heuristic:

A split is good if and only if: $SV_{i,g} > k * SV_{i,b}$, for good supervoxels $SV_{i,g}$ and bad supervoxels $SV_{i,b}$ for all relevant roots $i \in \{R_1, \dots, R_n\}$ larger than threshold t , within a local cutout of extent e . Relevant roots are any ground truth roots R_i who overlap the root that is split. For each of these, the split component that they overlap more is considered the base component, and the other component is considered the one that is split off. Good supervoxels are supervoxels of the split-off component that do not overlap with the root, bad supervoxels are those that do (and are erroneously split off). The precision factor k controls how many times more good supervoxels than bad supervoxels each root must have. These parameters were tuned based on authors' judgment to the following values: $k=20$, $t=1000$, $e=5000\text{nm}$.

Collecting error candidates

Using root skeletons fetched from the server, or generated via kimimaro, error candidates were generated heuristically: Skeleton nodes of degree 1 were treated as endpoints and split error candidates, skeleton nodes of degree > 2 were treated as junctions and merge error candidates.

To determine precision and recall of error candidates, actual errors and candidates were matched with [THRESHOLD] distance threshold.

Rendering

All samples were rendered using a unified mesh rendering pipeline relying on octarine; participating meshes were downloaded from CloudVolume. All renderings were centered on coordinates of errors and error candidates, with a view extent of 5000-8000nm.

Visualization of splits

To visualize splits, split components were matched to low-level meshes (L2 meshes) based on ChunkedGraph, which were rendered in distinct colors. Low-level meshes that were ambiguously associated with both components were colored by bisecting them with a plane orthogonal to the vector between the two sides' outer vertices that touched each colored mesh.

Tab 4

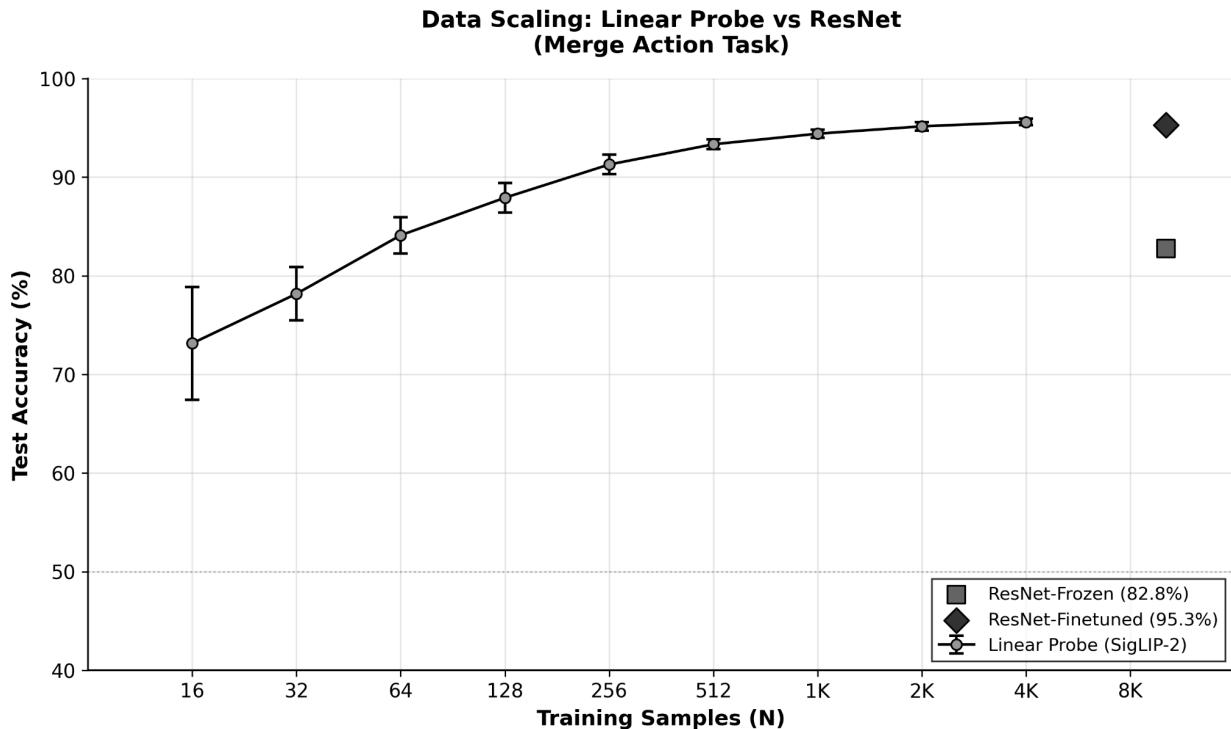
Opening framing: "We evaluate our approach along five dimensions: (1) how model scale affects data efficiency, (2) what capabilities emerge with full vision-language models, (3) whether language enables zero-shot adaptation, (4) whether learned priors generalize across species and imaging modalities, and (5) performance relative to human annotators. Together, these results demonstrate that vision-language models provide not just strong performance, but data-efficient, interpretable, adaptable, and generalizable proofreading capabilities."

1. How does model scale affect data efficiency? (*Core scaling insight*)

Goal: Establish that multimodal pretraining provides powerful visual priors

Show:

- Scaling curves (performance vs. # training examples)
 - ResNet: needs full dataset, struggles with few examples
 - SigLIP linear probe: strong performance with just 512 examples
 - Qwen VLM: best performance across all data regimes
- Comparison table at key data points (512, 2048, full dataset)
- Emphasize: frozen 400M param encoder with linear probe rivals ResNet trained on 10x more data



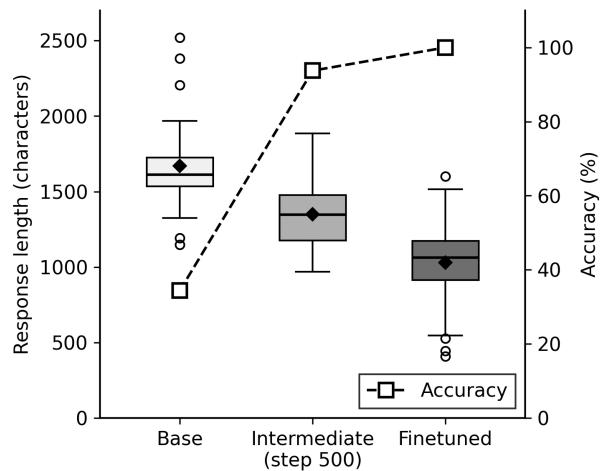
Message: "Multimodal pretraining provides powerful priors - even frozen encoders are highly data-efficient. Scale unlocks both performance and new capabilities."

2. What capabilities emerge with full VLM scale? (*Beyond classification*)

Goal: Show that full VLMs provide qualitatively different capabilities than smaller models

2.1 Reasoning: Chain-of-thought explanations

- Response style evolution during training (page 8 figure)
- Base model → intermediate → fine-tuned: response length changes
- Model learns to provide structured justifications

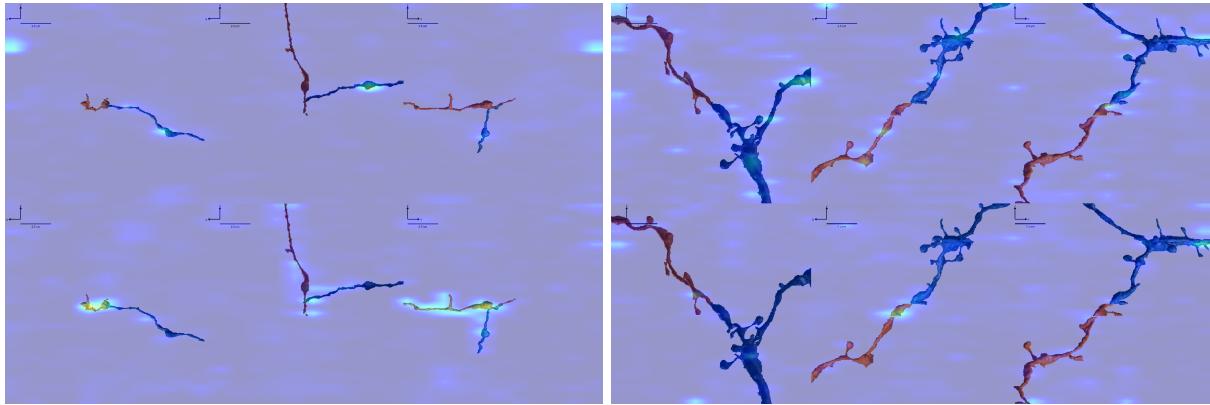


2.2 Interpretability: Understanding what the model learns

Visual attention (where):

- LayerCAM on SigLIP encoder + linear probe
- Show model attends to junction points, branch morphology, membrane boundaries
- Compare correct vs incorrect predictions

To generate class activation maps (CAMs) for the trained linear probe with the SigLIP vision encoder, we used the pytorch-grad-cam library, which computes gradients of class logits with respect to intermediate layer activations. Following a sweep over ViT layers and CAM methods (EigenCAM, GradCAM, LayerCAM, EigenGradCAM), we found LayerCAM on early-to-mid layers (Layers 1-6 out of 27) to elicit the most interpretable spatial attributions, specifically the layer normalization following the attention mechanism (layer_norm2) [cf FIGURE].



LayerCAM of ViT + linear probe (layer 2, left image: False example, right image: True example, Top is yes support, bottom is no support)

Semantic reasoning (why):

- Heuristics gained through training (page 8 analysis)
- Evolution from simple ("intact boundaries", "proximity without contact") to complex ("morphological semantics", "3D relational context", "depth/volumetric reasoning")
- Show example reasoning chains from base → intermediate → fine-tuned checkpoints
- What's gained: morphological semantics, multi-view alignment rigor, branch-pattern continuity
- What's lost: micro-boundary detection, projection-based hesitation, naive branch cues

Together these show:

- The model learns the "right" visual features (not shortcuts)
- And develops human-interpretable reasoning strategies
- Both improve with scale/training

2.3 Calibration: Confidence tracking accuracy

- Confidence distributions (page 7 plots)
- High confidence predictions more accurate than low confidence
- Model "knows when it doesn't know"
- Particularly important for positive samples (should merge) vs negative samples (should not merge)

Message: *"Full VLMs don't just classify better - they reason explainably, attend to interpretable features, and provide calibrated uncertainty estimates. These are emergent properties of scale."*

3. Can language enable zero-shot adaptation? (*Flexibility without retraining*)

Goal: Demonstrate unique advantage of language interface for incorporating expert knowledge

Show:

- Baseline: standard task prompt (used in training)
- Adaptive prompting experiments:
 - Conservative prompting: "Only answer 'yes' if highly confident"
 - Context injection: "This neuron is from [brain region/cell type]"
 - Heuristic guidance: "Pay special attention to [morphological feature]"
 - Few-shot examples in prompt
- Performance changes with different prompts (precision/recall trade-offs)
- Whether domain context improves accuracy on OOD data
- **Key comparison:** ResNet/linear probe require retraining; VLM adapts through language alone

Message: "*Language interface allows incorporating expert knowledge and adapting to new contexts without retraining - a unique advantage of VLMs over traditional CV.*"

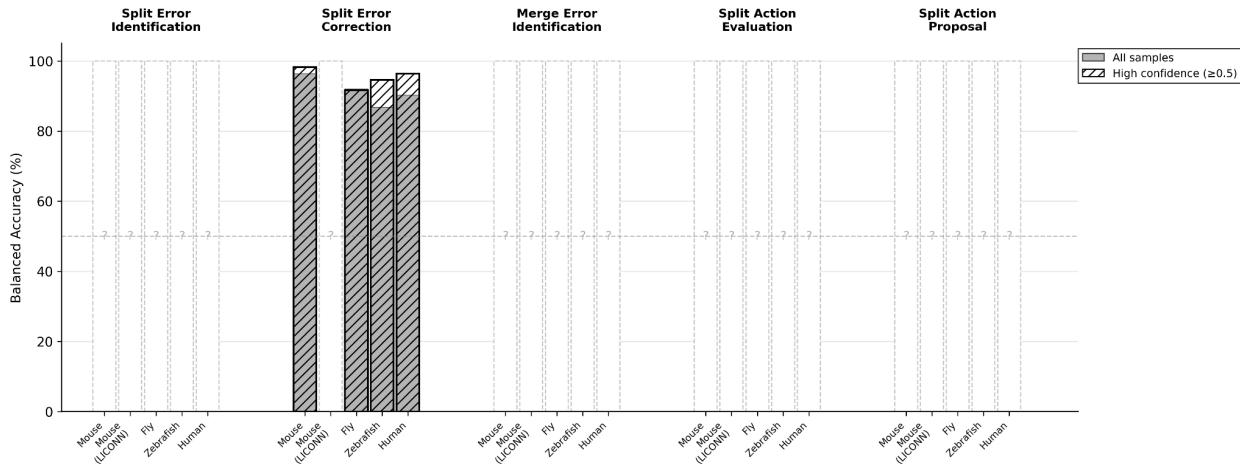
4. Do learned priors generalize? (*Cross-species/modality*)

Goal: Demonstrate that priors learned from mouse data transfer broadly

Show:

- Cross-species generalization:
 - Mouse (MICrONS) → Fly (FlyWire)
 - Mouse → Zebrafish
 - Mouse → Human (H01)
- Cross-modality generalization:
 - EM → ExM (expansion microscopy)
- Performance on each dataset (page 6-7 figures)
- Model trained **only on mouse** maintains strong performance
- Compare to baselines if available

Message: "*Priors from multimodal pretraining and scale generalize across biology and imaging modalities. The model learns fundamental visual reasoning about 3D neuronal structure, not dataset-specific patterns.*"



5. Human-level performance and beyond

Goal: Establish practical viability and competitive performance

Show:

- Within-distribution performance on MICrONS (page 6 figures)
- Comparison to human annotators (page 6 table):
 - Split Error Identification: $58.1\% \pm 6.3\%$
 - Split Error Correction: $91.9\% \pm 3.5\%$
 - Merge Error Identification: $72.6\% \pm 5.7\%$
 - Split Action Evaluation: $76.6\% \pm 5.3\%$
- Human inter-rater agreement
- GPU computational cost analysis (appendix)
- Integration feasibility into proofreading workflows

Message: "VLMs match or exceed human performance on several tasks, with feasible computational costs for practical deployment."

Closing synthesis: "Together, these results demonstrate a clear scaling trajectory: traditional CV methods are data-hungry and task-specific; frozen multimodal encoders provide data-efficient learning through strong visual priors; full vision-language models add reasoning, adaptability, and calibration. The result is a system that not only performs well but provides the interpretability, flexibility, and reliability needed for real-world scientific deployment."

Tab 6

Connectomics—mapping neural connectivity at nanometer resolution—requires extensive human proofreading to correct segmentation errors in electron and expansion microscopy data. We present ConnectomeVLM, a vision-language system that achieves human-level performance on proofreading tasks while demonstrating robust generalization and calibrated uncertainty. We fine-tune Qwen3-VL-32B using LoRA on multi-view orthographic renderings of neuronal segments, decomposing proofreading into four binary classification tasks: split error identification, split error correction, merge error identification, and merge error correction. On MICrONS mouse visual cortex data, our system achieves $91.9 \pm 3.5\%$ accuracy on split correction

and $72.6 \pm 5.7\%$ on merge identification, matching human annotator performance. Without retraining, the model generalizes across species (*Drosophila*, zebrafish, human) and imaging modalities (EM to expansion microscopy), maintaining $>85\%$ accuracy on out-of-distribution datasets. We demonstrate calibrated behavior through response variance analysis: the model's uncertainty correlates with dataset dissimilarity and predicts accuracy on negative samples. Analysis of learned heuristics reveals the model acquires morphologically-grounded reasoning—shifting from naive boundary detection to complex 3D spatial relationships, texture consistency, and branch pattern continuity—while maintaining interpretable, faithful descriptions of visual features. Our results establish vision-language models as

viable tools for scientific data curation tasks requiring expert-level visual reasoning.

Introduction

Understanding complex systems requires mapping their underlying structure. In neuroscience, this principle is especially acute: to understand how brains produce behavior, disease, and cognition, we need ground-truth maps of neural connectivity. Connectomics, the systematic mapping of neural connectivity at the level of individual synapses, is poised to provide these maps. By imaging entire nervous systems at nanometer resolution using electron microscopy or expansion microscopy, researchers can trace how each neuron connects to others, creating complete wiring diagrams of the brain. These high-resolution maps promise to transform neuroscience by enabling simulation of nervous systems, identification of structural signatures of disease, and reverse-engineering of the computational principles underlying biological intelligence.

However, a critical bottleneck limits progress: while automated segmentation algorithms can partition nanometer-resolution brain images into individual neurons, they make systematic errors that require extensive manual proofreading. These errors fall into two categories (Figure 1): *split errors*, where a single neuron is incorrectly fragmented into multiple segments, and *merge errors*, where parts of distinct neurons are incorrectly fused together. The FlyWire *Drosophila* connectome, for instance, required substantial human effort to correct 139,255 neurons, and scaling this approach to mammalian brains with tens of millions of neurons threatens to make whole-brain reconstruction economically infeasible.

Recent work has shown that large-scale vision-language models (VLMs) can perform proofreading tasks with zero-shot prompting. ConnectomeBench [cite] demonstrated that frontier models like o4-mini and Claude Sonnet 4 achieve competitive performance on tasks including identifying segmentation errors, validating proposed corrections, and determining whether neuron fragments should be merged. Surprisingly, these models solved 3D spatial reasoning problems using only 2D orthogonal projections, mirroring how human proofreaders visually inspect neurons and suggesting that pre-trained visual representations contain sufficient structure to support connectomics workflows. However, these frontier models are large (>100B parameters), proprietary, and expensive to deploy at scale. This raises fundamental questions about what computational resources are actually necessary: What function does language play

versus pure vision? What role do pre-trained representations play versus task-specific fine-tuning? What is the contribution of model capacity versus training data scale? Are frontier-scale models strictly necessary, or can smaller, task-specialized models achieve comparable performance?

This paper provides a systematic scaling analysis to answer these questions. We compare three model classes across four proofreading subtasks that span a range of reasoning complexity: (1) vision-only models (CNNs and ViTs), (2) vision-language models (VLMs) without generation (linear probes on SigLip), and (3) generative VLMs (fine-tuned with LoRA). We evaluate along four deployment-critical axes: performance (data requirements to achieve human-level accuracy), generalization (zero-shot transfer to new species and imaging modalities), calibration (reliability of uncertainty estimates for human-AI collaboration), and interpretability (whether models learn meaningful heuristics). This analysis reveals task-specific scaling laws that inform the design of economically viable automated proofreading systems.

Our systematic evaluation reveals task-specific scaling laws: simple topology-matching tasks achieve human-level performance with vision-only models trained on fewer than 1,000 samples, while complex reasoning tasks benefit from larger vision-language models but still achieve strong performance with task-specific fine-tuning. We demonstrate robust cross-species generalization, with models trained on mouse cortex transferring to *Drosophila*, zebrafish, and human datasets with modest accuracy degradation. We show that models can be well-calibrated, enabling reliable uncertainty quantification for human-AI collaboration, and that they learn biologically meaningful visual heuristics rather than exploiting spurious shortcuts. Finally, we demonstrate that an integrated system combining task-appropriate models can perform end-to-end proofreading at a computational cost of \$600 to \$2,000 per dataset (orders of magnitude lower than human annotation), making AI-assisted whole-brain reconstruction economically viable.

The remainder of the paper is organized as follows. Section 2 describes our training methods and data generation pipeline. Section 3 presents results organized by evaluation axis: performance and scaling behavior, cross-species generalization, calibration analysis, and interpretability findings. Section 4 demonstrates an integrated proofreading system and estimates computational costs for large-scale deployment. Section 5 discusses implications for connectomics workflows and future AI-assisted scientific annotation tasks.

Proofreading Tasks and Difficulty Hierarchy

Connectomics segmentation errors fall into two categories: split errors, where a single neuron is fragmented into multiple segments, and merge errors, where parts of multiple neurons are incorrectly fused. The proofreading workflow decomposes into three stages: (1) error candidate localization (identifying potential error sites using skeleton-based heuristics), (2) error identification (determining whether a candidate represents a true error), and (3) error correction (validating proposed fixes). Our skeleton-based heuristics achieve approximately 95% recall for

split errors and 60% recall for merge errors, though at only 2% precision, necessitating accurate identification models.

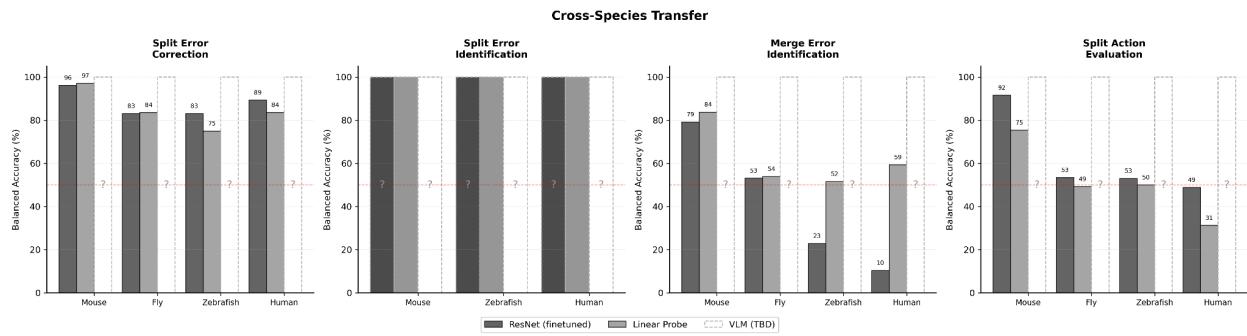
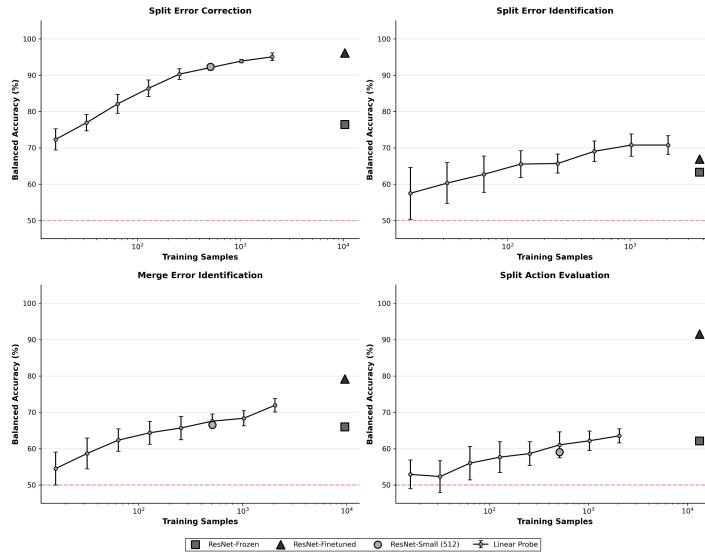
We focus on identification and correction tasks across a difficulty spectrum:

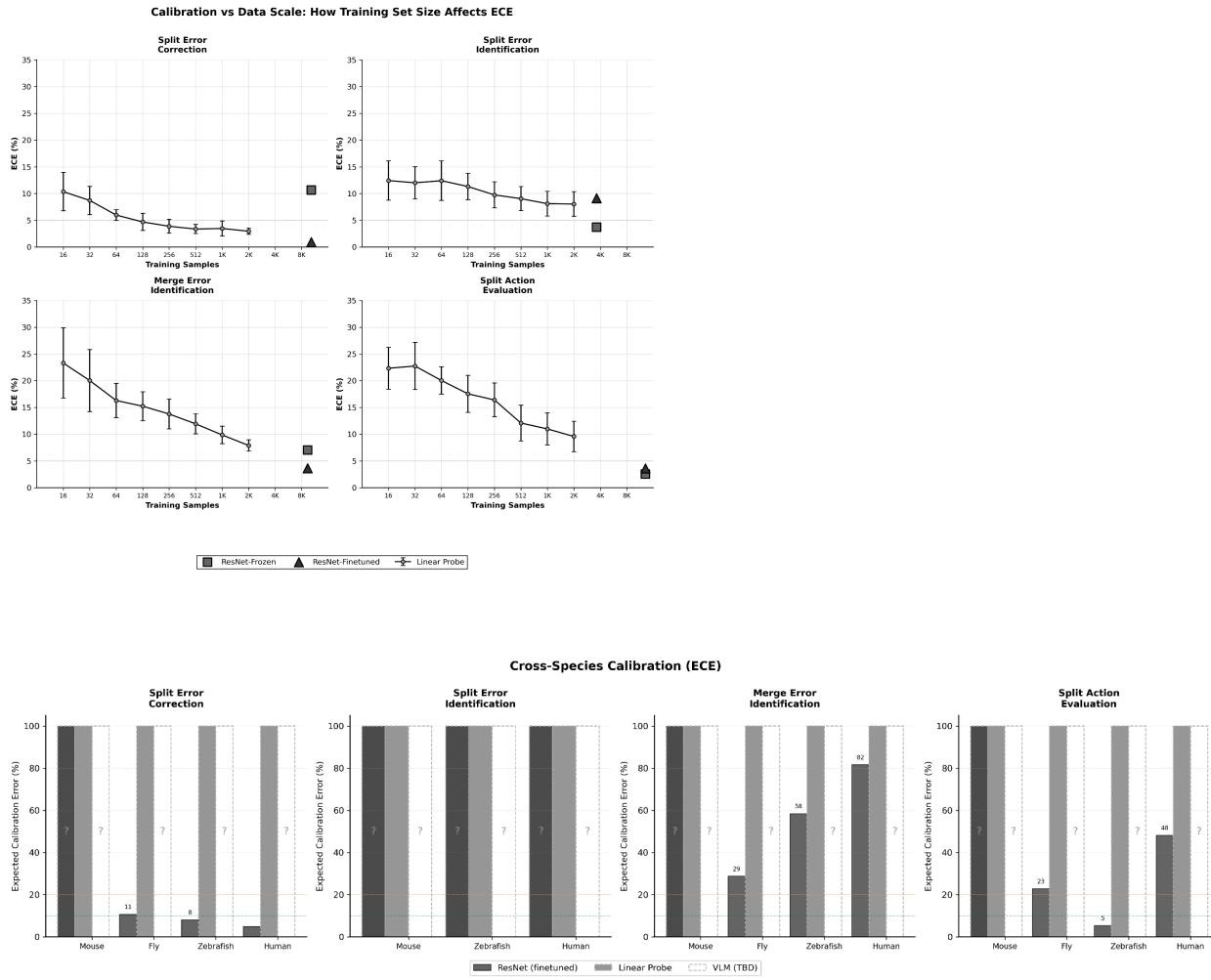
Split Error Correction (Low complexity): Given two segment endpoints, determine whether they should be merged. Requires local topology matching by evaluating whether segments exhibit geometric continuity and consistent morphology at the junction point. Success depends on recognizing alignment across multiple 2D views and matching branch patterns.

Split Error Identification (Medium complexity): Given a segment endpoint, determine whether it represents a true split error requiring correction. Requires multimodal reasoning, fusing information from both 3D segmentation renderings and raw electron microscopy (EM) or expansion microscopy (ExM) images. The model must detect whether the endpoint corresponds to a true neuronal terminus or an artificial break introduced by segmentation errors.

Merge Error Correction Tasks (High complexity):

- **Merge Error Identification:** Given a segment, determine whether it contains a merge error where parts of multiple distinct neurons have been incorrectly fused. Requires global visual reasoning over the entire segment morphology to detect discontinuities, inconsistent branching patterns, or regions where disparate neuronal structures have been incorrectly joined.
- **Split Action Evaluation:** Given a proposed split correction, determine whether it successfully separates merged neurons. Requires counterfactual visual reasoning. The model must evaluate whether the split-off component genuinely belongs to a different neuron, often by assessing whether morphological patterns diverge across the proposed boundary.





Introduction

Start with something that everyone agrees on | Connectomics is the systematic mapping of neural connectivity at the level of individual synapses. By imaging entire nervous systems at nanometer resolution using electron microscopy or expansion microscopy, researchers can trace how each neuron connects to others, creating complete wiring diagrams of the brain. These high-resolution maps promise to transform neuroscience by enabling simulation of nervous systems, identification of structural signatures of disease, and reverse-engineering of

the computational principles underlying biological intelligence. However, a critical bottleneck limits progress: while automated segmentation algorithms can partition nanometer-resolution brain images into individual neurons, they make systematic errors [Introduce an example](#) that require extensive manual proofreading. The FlyWire *Drosophila* connectome, for instance, required substantial human effort to correct 139,255 neurons, and scaling this approach to mammalian brains with tens of millions of neurons threatens to make whole-brain reconstruction economically infeasible.

Recent work has shown that large-scale vision-language models (VLMs) can perform proofreading tasks with zero-shot prompting. ConnectomeBench [cite] demonstrated that frontier models like o4-mini and Claude Sonnet 4 achieve competitive performance on tasks including identifying segmentation errors, validating proposed corrections, and determining whether neuron fragments should be merged. Surprisingly, these models solved 3D spatial reasoning problems using only 2D orthogonal projections, suggesting that pre-trained visual and linguistic representations contain sufficient structure to support connectomics workflows. However, these frontier models are large (>100B parameters), proprietary, and expensive to deploy at scale. This raises fundamental questions about what computational resources are actually necessary: What function does language play versus pure vision? What role do pre-trained representations play versus task-specific fine-tuning? What is the contribution of model capacity versus training data scale? Are frontier-scale models strictly necessary, or can smaller, task-specialized models achieve comparable performance?

This paper provides a systematic scaling analysis to answer these questions. We compare three model classes across four proofreading subtasks that span a range of reasoning complexity: **(1) vision-only models** (CNNs and ViTs), **(2) vision-language models without generation** (linear probes on SigLip), and **(3) generative VLMs** (fine-tuned with LoRA, capable of producing natural language explanations). We evaluate along four deployment-critical axes: performance (data requirements to achieve human-level accuracy), generalization (zero-shot transfer to new species and imaging modalities), calibration (reliability of uncertainty estimates for human-AI collaboration), and interpretability (whether models learn meaningful heuristics). This analysis reveals task-specific scaling laws that inform the design of economically viable automated proofreading systems.

Proofreading Tasks and Difficulty Hierarchy

Connectomics segmentation errors fall into two categories: **split errors**, where a single neuron is fragmented into multiple segments, and **merge errors**, where parts of multiple neurons are incorrectly fused. The proofreading workflow decomposes into three stages: (1) error candidate localization (identifying potential error sites using skeleton-based heuristics), (2) error identification (determining whether a candidate represents a true error), and (3) error correction (validating proposed fixes). Our skeleton-based heuristics achieve approximately 95% recall for split errors and 60% recall for merge errors, though at only 2% precision, necessitating accurate identification models.

We focus on identification and correction tasks across a difficulty spectrum:

Split Error Correction (Low complexity): Given two segment endpoints, determine whether they should be merged. Requires local topology matching by evaluating whether segments exhibit geometric continuity and consistent morphology at the junction point. Success depends on recognizing alignment across multiple 2D views and matching branch patterns.

Split Error Identification (Medium complexity): Given a segment endpoint, determine whether it represents a true split error requiring correction. Requires multimodal reasoning, fusing information from both 3D segmentation renderings and raw electron microscopy (EM) or expansion microscopy (ExM) texture. The model must detect whether the endpoint corresponds to a true neuronal terminus or an artificial break introduced by segmentation errors.

Merge Error Correction Tasks (High complexity):

- **Merge Error Identification:** Given a segment, determine whether it contains a merge error where parts of multiple distinct neurons have been incorrectly fused. Requires global visual reasoning over the entire segment morphology to detect discontinuities, inconsistent branching patterns, or regions where disparate neuronal structures have been incorrectly joined.
- **Split Action Evaluation:** Given a proposed split correction, determine whether it successfully separates merged neurons. Requires counterfactual visual reasoning—the model must evaluate whether the split-off component genuinely belongs to a different neuron, often by assessing whether morphological patterns diverge across the proposed boundary.

Why Image-Based Models?

We adopt an image-based approach for three reasons. First, human proofreaders perform these tasks through visual inspection of 2D projections, suggesting that image-based reasoning is sufficient for the problem. Second, pre-trained vision models have developed rich visual priors that generalize across domains, and we aim to leverage these representations rather than engineer task-specific graph features. Third, vision-language models benefit from favorable scaling laws, allowing us to study how linguistic grounding affects performance on fundamentally visual tasks. While graph-based methods (operating on meshes or skeletons) are complementary, they do not directly benefit from the pre-training and scaling properties of large vision-language models.

Our Approach and Contributions

We systematically evaluate how three factors—model architecture, training data scale, and task complexity—interact to determine proofreading performance. Our contributions are:

Task-specific scaling laws: We measure how data requirements vary across proofreading subtasks and identify when simple vision-only models suffice versus when language and generation capabilities provide benefits.

Cross-species generalization: We quantify transfer performance from mouse cortex to other species (*Drosophila*, zebrafish) and imaging modalities (expansion microscopy), demonstrating that models trained on one organism generalize to others with modest accuracy degradation (X-Y% depending on task complexity).

Calibration analysis: We assess whether models can reliably communicate uncertainty, which is critical for human-AI collaboration. We find that [placeholder for key calibration finding].

Interpretability: Using gradient-based attribution and analysis of VLM reasoning chains, we show that models learn biologically meaningful heuristics (morphological continuity, trajectory independence) rather than exploiting spurious shortcuts. We identify which heuristics are gained and lost during training.

Economic viability: We demonstrate that an integrated system combining task-appropriate models can perform end-to-end proofreading at a computational cost of \$600-\$2,000 per dataset, orders of magnitude lower than human annotation costs, making AI-assisted whole-brain reconstruction economically feasible.

The remainder of the paper is organized as follows. Section 2 describes our training methods and data generation pipeline. Section 3 presents results organized by evaluation axis: performance and scaling behavior, cross-species generalization, calibration analysis, and interpretability findings. Section 4 demonstrates an integrated proofreading system and estimates computational costs for large-scale deployment. Section 5 discusses implications for connectomics workflows and future AI-assisted scientific annotation tasks.

Introduction

Connectomics is the systematic mapping of neural connectivity at the level of individual synapses. By imaging entire nervous systems at nanometer resolution using electron

microscopy or expansion microscopy, researchers can trace how each neuron connects to others, creating complete wiring diagrams of the brain. These high-resolution maps promise to transform neuroscience by enabling simulation of nervous systems, identification of structural signatures of disease, and reverse-engineering of the computational principles underlying biological intelligence. However, a critical bottleneck limits progress: while automated segmentation algorithms can partition nanometer-resolution brain images into individual neurons, they make systematic errors that require extensive manual proofreading. The FlyWire *Drosophila* connectome, for instance, required substantial human effort to correct 139,255 neurons, and scaling this approach to mammalian brains with tens of millions of neurons threatens to make whole-brain reconstruction economically infeasible.

Recent work has shown that large-scale vision-language models (VLMs) can perform proofreading tasks with zero-shot prompting. ConnectomeBench [cite] demonstrated that frontier models like o4-mini and Claude Sonnet 4 achieve competitive performance on tasks including identifying segmentation errors, validating proposed corrections, and determining whether neuron fragments should be merged. Surprisingly, these models solved 3D spatial reasoning problems using only 2D orthogonal projections, suggesting that pre-trained visual and linguistic representations contain sufficient structure to support connectomics workflows. However, these frontier models are large (>100B parameters), proprietary, and expensive to deploy at scale. This raises fundamental questions about what computational resources are actually necessary: What function does language play versus pure vision? What role do pre-trained representations play versus task-specific fine-tuning? What is the contribution of model capacity versus training data scale? Are frontier-scale models strictly necessary, or can smaller, task-specialized models achieve comparable performance?

This paper provides a systematic scaling analysis to answer these questions. We compare three model classes across four proofreading subtasks that span a range of reasoning complexity: **(1) vision-only models** (CNNs, ViTs, linear probes on frozen vision encoders like ResNet and SigLIP), **(2) vision-language models without generation** (linear probes on contrastively pre-trained encoders that were trained with image-text pairs but perform classification without text generation), and **(3) generative VLMs** (fine-tuned with LoRA, capable of producing natural language explanations). We evaluate along four deployment-critical axes: performance (data requirements to achieve human-level accuracy), generalization (zero-shot transfer to new species and imaging modalities), calibration (reliability of uncertainty estimates for human-AI collaboration), and interpretability (whether models learn biologically meaningful heuristics). This analysis reveals task-specific scaling laws that inform the design of economically viable automated proofreading systems.

Proofreading Tasks and Difficulty Hierarchy

Connectomics segmentation errors fall into two categories: **split errors**, where a single neuron is fragmented into multiple segments, and **merge errors**, where parts of multiple neurons are incorrectly fused. The proofreading workflow decomposes into three stages: (1) error candidate localization (identifying potential error sites using skeleton-based heuristics), (2) error

identification (determining whether a candidate represents a true error), and (3) error correction (validating proposed fixes).

We focus on identification and correction tasks across a difficulty spectrum:

Split Error Correction (Low complexity): Given two segment endpoints, determine whether they should be merged. Requires local topology matching by evaluating whether segments exhibit geometric continuity and consistent morphology at the junction point. Success depends on recognizing alignment across multiple 2D views and matching branch patterns.

Split Error Identification (Medium complexity): Given a segment endpoint, determine whether it represents a true split error requiring correction. Requires multimodal reasoning, fusing information from both 3D segmentation renderings and raw electron microscopy (EM) or expansion microscopy (ExM) texture. The model must detect whether the endpoint corresponds to a true neuronal terminus or an artificial break introduced by segmentation errors.

Merge Error Correction Tasks (High complexity):

- **Merge Error Identification:** Given a segment, determine whether it contains a merge error where parts of multiple distinct neurons have been incorrectly fused. Requires global visual reasoning over the entire segment morphology to detect discontinuities, inconsistent branching patterns, or regions where disparate neuronal structures have been incorrectly joined.
- **Split Action Evaluation:** Given a proposed split correction, determine whether it successfully separates merged neurons. Requires counterfactual visual reasoning—the model must evaluate whether the split-off component genuinely belongs to a different neuron, often by assessing whether morphological patterns diverge across the proposed boundary.

Why Image-Based Models?

We adopt an image-based approach for three reasons. First, human proofreaders perform these tasks through visual inspection of 2D projections, suggesting that image-based reasoning is sufficient for the problem. Second, pre-trained vision models have developed rich visual priors that generalize across domains, and we aim to leverage these representations rather than engineer task-specific graph features. Third, vision-language models benefit from favorable scaling laws, allowing us to study how linguistic grounding affects performance on fundamentally visual tasks. While graph-based methods (operating on meshes or skeletons) are complementary, they do not directly benefit from the pre-training and scaling properties of large vision-language models.

Our Approach and Contributions

We systematically evaluate how three factors—model architecture, training data scale, and task complexity—interact to determine proofreading performance. Our contributions are:

Task-specific scaling laws: We measure how data requirements vary across proofreading subtasks and identify when simple vision-only models suffice versus when language and generation capabilities provide benefits.

Cross-species generalization: We quantify transfer performance from mouse cortex to other species (*Drosophila*, zebrafish) and imaging modalities (expansion microscopy), demonstrating that models trained on one organism generalize to others with modest accuracy degradation (X-Y% depending on task complexity).

Calibration analysis: We assess whether models can reliably communicate uncertainty, which is critical for human-AI collaboration. We find that [placeholder for key calibration finding].

Interpretability: Using gradient-based attribution and analysis of VLM reasoning chains, we show that models learn biologically meaningful heuristics (morphological continuity, trajectory independence) rather than exploiting spurious shortcuts. We identify which heuristics are gained and lost during training.

Economic viability: We demonstrate that an integrated system combining task-appropriate models can perform end-to-end proofreading at a computational cost of \$600-\$2,000 per dataset, orders of magnitude lower than human annotation costs, making AI-assisted whole-brain reconstruction economically feasible.

The remainder of the paper is organized as follows. Section 2 describes our training methods and data generation pipeline. Section 3 presents results organized by evaluation axis: performance and scaling behavior, cross-species generalization, calibration analysis, and interpretability findings. Section 4 demonstrates an integrated proofreading system and estimates computational costs for large-scale deployment. Section 5 discusses implications for connectomics workflows and future AI-assisted scientific annotation tasks.

Introduction

Stop using all these em-dashes — lol

Connectomics—the systematic mapping of neural connectivity—promises to transform neuroscience by enabling simulation of nervous systems, identification of disease biomarkers, and reverse-engineering of biological intelligence. We need more context here. The reviewer probably hasn't heard of connectomics in their life However, a critical bottleneck limits progress: while automated segmentation algorithms can partition nanometer-resolution brain images into individual neurons, they make systematic errors that require extensive manual proofreading. The FlyWire *Drosophila* connectome, for instance, required substantial human effort to correct 139,255 neurons, and scaling this approach to mammalian brains with tens of millions of neurons threatens to make whole-brain reconstruction economically infeasible.

Recent work has shown that large-scale vision-language models (VLMs) can perform proofreading tasks with zero-shot prompting. ConnectomeBench [cite] demonstrated that frontier models like o4-mini and Claude Sonnet 4 achieve competitive performance on various proofreading subtasks Confusing to say here without mentioning the subtasks tasks, maybe can quickly describe or find another way to communicate, surprisingly solving 3D spatial reasoning problems using only 2D orthogonal projections. These results suggest that pre-trained visual and linguistic representations contain sufficient structure to support connectomics workflows. However, critical questions remain unanswered: Are frontier-scale proprietary models strictly necessary, or can smaller, task-specialized models achieve comparable performance? The transition here is unclear, how does this question from. Probably benefits for better set up in teh previous sentence. Another related questions that we can perhaps synthesize: What function is language playing? What functions are the representations from pretraining playing? What function is model capacity playing? What are the tradeoffs between vision-only, vision-language, and generative architectures? How do data requirements, generalization capabilities, and uncertainty quantification vary across different proofreading subtasks?

This paper provides a systematic scaling analysis to answer these questions. We compare three model classes—vision-only models (CNNs, ViTs, linear probes on frozen encoders), vision-language models without generation (linear probes on contrastively pre-trained encoders), and generative VLMs (fine-tuned with LoRA)—across four proofreading subtasks that span a range of reasoning complexity. We evaluate along four deployment-critical axes: **(1) performance** (data requirements to achieve human-level accuracy), **(2) generalization** (zero-shot transfer to new species and imaging modalities), **(3) calibration** (reliability of uncertainty estimates), and **(4) interpretability** (whether models learn biologically meaningful heuristics). This analysis reveals task-specific scaling laws that inform the design of economically viable automated proofreading systems.

Proofreading Tasks and Difficulty Hierarchy

Connectomics segmentation errors fall into two categories: **split errors**, where a single neuron is fragmented into multiple segments, and **merge errors**, where parts of multiple neurons are incorrectly fused. The proofreading workflow decomposes into three stages: (1) **error candidate**

localization, (2) **error identification** (determining whether a candidate represents a true error), and (3) **error correction** (validating proposed fixes).

We focus on identification and correction tasks across a difficulty spectrum:

- **Split Error Correction** (Low complexity): Given two segment endpoints, determine whether they should be merged. Requires local topology matching—evaluating whether segments exhibit geometric continuity and consistent morphology at the junction point. Success depends on recognizing alignment across multiple 2D views and matching branch patterns.
- **Split Error Identification** (Medium complexity): Given a segment endpoint, determine whether it represents a true split error requiring correction. Requires multimodal reasoning, fusing information from both 3D segmentation renderings and raw electron microscopy (EM) or expansion microscopy (ExM) texture. The model must detect whether the endpoint corresponds to a true neuronal terminus or an artificial break introduced by segmentation errors.
- **Merge Error Correction Related** (High complexity)
 - **Merge Error Identification**: Given a segment, determine whether it contains a merge error where parts of multiple distinct neurons have been incorrectly fused. Requires global visual reasoning over the entire segment morphology to detect discontinuities, inconsistent branching patterns, or regions where disparate neuronal structures have been incorrectly joined.
 - **Split Action Evaluation** Given a proposed split correction, determine whether it successfully separates merged neurons. Requires counterfactual visual reasoning—the model must evaluate whether the split-off component genuinely belongs to a different neuron, often by assessing whether morphological patterns diverge across the proposed boundary.

Evaluation Framework: Four Axes for Deployment Readiness

Placement here is weird

Performance: We measure class-balanced accuracy and compare against human annotator baselines. The central question is how much training data each model class requires to achieve human-level performance on each task. This determines the annotation cost and feasibility of creating task-specific training sets.

Generalization: We train models on mouse cortex data (MICrONS dataset) and evaluate zero-shot transfer to other species (*Drosophila* from FlyWire, zebrafish) and imaging modalities (expansion microscopy from LICONN). Strong generalization is essential for practical deployment, as re-annotating large training sets for each new dataset is prohibitively expensive.

Calibration: We compute Expected Calibration Error (ECE) and Maximum Calibration Error (MCE) to assess whether model confidence scores reliably indicate prediction uncertainty. Calibration is critical for human-AI collaboration: well-calibrated models can flag uncertain predictions for human review, preventing the introduction of new errors that would require additional proofreading. For vision-only models, we evaluate calibration using output logits; for generative VLMs, we assess verbalized confidence estimates.

Interpretability and Steering: We use gradient-based attribution methods (GradCAM) to visualize which image regions influence predictions, and we analyze the reasoning patterns that emerge in generative VLMs through qualitative examination of generated explanations. This reveals whether models learn biologically meaningful heuristics (e.g., morphological continuity, trajectory independence) or exploit spurious shortcuts. Interpretability also enables steering—identifying and correcting failure modes without full retraining.

Key Findings

Our analysis reveals task-dependent scaling laws and architecture-specific tradeoffs. For simple topology-matching tasks, lightweight vision-only models achieve near-human performance with minimal training data. For multimodal fusion and complex reasoning tasks, larger vision-language models provide substantial benefits, but generative capabilities are not always necessary. We demonstrate strong cross-species generalization, with models trained on mouse data transferring to other organisms with modest performance degradation. **Maybe this happens**
Finally, we show that an integrated system combining task-appropriate models can perform end-to-end proofreading at a computational cost orders of magnitude lower than human annotation, making AI-assisted whole-brain reconstruction economically viable.

The remainder of the paper is organized as follows. Section 2 describes our training methods and data generation pipeline. Section 3 presents results organized by evaluation axis: performance and scaling behavior, cross-species generalization, calibration analysis, and interpretability findings. Section 4 demonstrates an integrated proofreading system and estimates computational costs for large-scale deployment. Section 5 discusses implications for connectomics workflows and future AI-assisted scientific annotation tasks.