

1 GPU Computational Cost Estimation for AI-Based Proofreading

1.1 Motivation and Methods

Connectome proofreading represents a major practical bottleneck in connectomics. FlyWire required 30 human-years to fully proofread the *Drosophila* brain’s 139,255 neurons; such timescales prohibit rapid iteration on segmentation algorithms and limit deployment to new datasets. To enable efficient AI-based proofreading at scale, we must understand the computational cost landscape. We analyzed edit histories from two major connectomics datasets: the MICrONS mouse cortex (2,314 proofreading-accessible neurons in a 1 mm³ volume, representing a partially-proofread mammalian circuit), and FlyWire’s *Drosophila* brain (139,255 neurons, fully proofread, representing a large-scale completed dataset).

Retrieving full edit histories for entire connectome datasets is computationally intensive, so we sampled strategically from the proofread subsets. We sampled proofread neurons from each dataset—mouse: $n = 500$ from 2,314 proofread neurons in the 1 mm³ MICrONS volume, fly: $n = 1,000$ from 139,255 proofread neurons in the complete FlyWire brain—and retrieved complete edit histories. For both species, differences between metrics on number of edits per neuron between $n = 500$ and $n = 1,000$ were small ($< 7\%$) providing confidence that linear extrapolation from these sample sizes to full proofread populations is statistically robust (Figure 1.4).

For each sampled neuron, we categorized operations as merge (consolidating over-split fragments) or split (separating under-segmented regions), computed edit count distributions, and identified heavy-tail concentration (95th percentile threshold). We then modeled GPU computational costs using two approaches:

Naive Model (uniform times): Assumes both merge and split operations require equal inference time:

$$\text{GPU Cost (Naive)} = \frac{(\text{Merge} + \text{Split}) \times t_{\text{uniform}}}{3600 \text{ sec/hr}} \times \$2/\text{GPU-hour} \quad (1)$$

where $t_{\text{uniform}} = 2.0$ seconds per operation.

Realistic Model (task-stratified times): Distinguishes merge and split complexity based on proofreading task difficulty (Section 3.1): merge corrections (consolidating multiple fragments) represent higher-complexity operations requiring more inference passes, while split corrections (isolating over-merged regions) require less:

$$\text{GPU Cost (Realistic)} = \frac{(\text{Merge} \times t_{\text{merge}} + \text{Split} \times t_{\text{split}})}{3600 \text{ sec/hr}} \times \$2/\text{GPU-hour} \quad (2)$$

where $t_{\text{merge}} = 2.5$ seconds (high complexity), $t_{\text{split}} = 1.5$ seconds (medium complexity), and total per-operation inference time ranges 1.5–2.5 seconds depending on operation distribution. Hardware: Qwen-32B model on dual H100 GPUs, GPU rate \$2/hour.

1.2 Computational Cost Estimation at Three Scales

We estimate proofreading costs at three extrapolation levels, progressing from current data to hypothetical scenarios (Table 1):

Level	Mouse	Cost (Naive / Realistic)	Fly	Cost (Naive / Realistic)
Level 1: Current proofread	2,314 neurons	\$1,189 / \$1,189 (594 / 594 GPU-hrs)	139,255 neurons	\$3,046 / (1,523 / 1,523 GPU-hrs)
Level 2: Connectome volume	~75,000 neurons	\$34,262 / \$33,628 (17,131 / 16,814 GPU-hrs)	139,255 neurons	\$3,046 / (1,364 / 1,528 GPU-hrs)
Level 3: Whole cortex/brain	~10,000,000 neurons	\$4,568,222 / \$4,483,710 (2,284,111 / 2,241,855 GPU-hrs)	~140,000 neurons	\$2,728 / (1,523 / 1,523 GPU-hrs)

Level 1 = current fully-analyzed proofread neurons; Level 2 = expected connectomic volume; Level 3 = hypothetical

Table 1: **GPU Computational Cost Estimation: Three Extrapolation Scales \times Two Models.** Naive model assumes uniform 2.0s per operation; Realistic model uses task-stratified times (merge=2.5s, split=1.5s). Mouse Level 2 assumes 75,000 expected neurons in 1 mm³ based on MICrONS density; Fly Level 2-3 are identical since brain is fully proofread. Level 3 for mouse illustrates 4.3 \times cost scaling for full mammalian cortex (112 mm³, ~10M neurons). GPU cost = \$2/hour.

The three levels reveal markedly different scaling profiles: mammalian proofreading costs scale dramatically with volume (4,300 \times from Level 1 to Level 3), while insect costs remain constant (fly brain is uniformly proofread). Our subsequent analyses focus on **Level 2 (connectome volume scale)** because it represents expected computational loads for typical connectomics deployments, balancing current data with realistic volume expectations.

1.3 Results

Edit distributions reveal species-level segmentation differences. Mouse neurons require dramatically more proofreading: 411 ± 288 edits per neuron (median=335) versus 17.5 ± 32 edits per neuron (median=8) for fly—a 23.4 \times intensity difference. This reflects distinct anatomical challenges: mammalian cortex has higher neuronal density and more complex morphology, leading to more substantial segmentation errors requiring correction (Figure 1A).

Operation ratios expose segmentation biases. Mouse proofreading exhibits balanced merge-to-split ratios (46.3% merge, 53.7% split), indicating that the initial segmentation contains comparable over- and under-segmentation errors. Fly proofreading is merge-dominated (74.0% merge, 26.0% split), revealing systematic oversegmentation in the FAFB volume where numerous small fragments must be consolidated. These divergent patterns suggest that cost-effective proofreading systems must handle both error modes (Figure 1B).

Heavy-tail distribution concentrates computational effort. Approximately 5% of neurons exceed the 95th percentile threshold, but their contribution differs markedly. Mouse heavy-tail neurons (>971 edits) account for only

14.1% of total edits, indicating relatively uniform workload distribution. In contrast, fly heavy-tail neurons (60 edits) concentrate 33.3% of total edits despite $23.4\times$ lower per-neuron baseline, revealing that fly proofreading is dominated by a small number of structurally complex outliers (Figure 1A).

Cost projections: Connectomic volume scale (Level 2 in Table 1).

Figure 1D shows GPU costs at the connectomic volume scale—expected neuron populations for deployments on each dataset—using the naive model (uniform 2.0 s per operation) for simplicity. This represents typical computational loads for full-volume proofreading systems:

- **Mouse (expected 75,000 neurons in 1 mm³ MICrONS volume):**
30,835,500 total edits \rightarrow 16,814 GPU-hours, \$33,628 cost
- **Fly (expected 140,000 neurons in connectomic volume):** 2,455,600 total edits \rightarrow 1,528 GPU-hours, \$3,056 cost

Despite $23.4\times$ lower per-neuron effort (17.5 vs 411 edits/neuron), fly costs only 4.6% less than mouse (\$3.1K vs \$33.6K) because mouse’s high per-neuron complexity is multiplied by a proportionally larger expected neuron population (75K vs 140K). This demonstrates that connectomic deployment costs are dominated by dataset size and segmentation quality rather than intrinsic anatomical complexity. The dramatic scaling becomes apparent at full-cortex scale: if the entire mouse neocortex (112 mm³, \sim 10 million neurons) required the same per-neuron proofreading effort, costs would scale to \$4.5 million (2.2 million GPU-hours, Level 3), underscoring why connectome proofreading at full mammalian brain scales remains computationally and economically prohibitive.

Cost sensitivity and realistic ranges. At connectomic volume scale, costs are highly sensitive to per-operation inference time. The realistic 1.5–2.5 second range (Figure 1C, highlighted by dashed rectangle) demonstrates that inference latency is the dominant cost driver. Within this realistic window, mouse costs span roughly \$19K–\$38K while fly costs span \$3K–\$6K, depending on inference time assumptions. Improving inference speed from 2.5 to 1.5 seconds saves 40% on GPU hours, illustrating why model optimization and algorithmic efficiency are critical for scalability at connectomic volumes.

Statistical validation. Cross-sample consistency between $n = 100$ and $n = 500$ mouse samples validates our extrapolation approach: mean edits differ by 2.5%, projected totals differ by 2.5%, and heavy-tail contribution differs by 32.6%. Similar consistency for fly ($n = 100$ vs $n = 1,000$: 6.6% and 6.5% differences) demonstrates robust sampling methodology.

1.4 Discussion

Our three-level cost analysis (Table 1) reveals a fundamental principle: GPU-accelerated proofreading costs scale with dataset volume, not species complexity. At Level 1 (current proofread data), mouse is cheaper (\$1,189 vs \$3,046). At Level 2 (connectomic volumes), mouse becomes $11\times$ more expensive (\$33,628 vs \$3,056) despite lower per-neuron effort, purely due to expected population size.

At Level 3 (full mouse cortex), costs exceed mammalian budgets (\$4.5M), while fly remains constant. This inverse relationship (simpler anatomy, higher cost) underscores why connectome proofreading at scale is economically prohibitive for mammals but feasible for insects.

The realistic computational model (task-stratified inference times) provides more accurate cost predictions than naive uniform models. By accounting for merge operations ($t_{\text{merge}} = 2.5$ s, high complexity) vs. split operations ($t_{\text{split}} = 1.5$ s, medium complexity)—distinctions grounded in proofreading task structure (Section 3.1)—we capture the true cost landscape. Mouse shows merge-dominated GPU burden (59%), while fly exhibits even stronger merge dominance (89%), reflecting systematic segmentation biases in FAFB. These patterns indicate that optimal proofreading systems should be adaptive: detecting error-type dominance and prioritizing computational resources accordingly.

The divergent merge/split error patterns also suggest that improving initial segmentation quality provides asymmetric payoffs. Reducing merge-heavy error types (as in fly) offers greater cost savings than reducing split-heavy errors (as in mouse), because merge corrections require more inference. The computationally feasible costs for current connectomics deployments (\$3K–\$34K for typical volumes) make GPU-accelerated proofreading economically viable compared to human annotation (30 human-years for FlyWire). However, scaling to larger mammalian brain volumes (full neocortex, zebrafish, primate) will require two complementary approaches: (1) algorithmic advances to reduce per-operation latency (e.g., specialized models for merge vs. split, batch processing), and (2) improved initial segmentation to reduce edit counts (e.g., advanced 3D U-Nets, attention mechanisms).

Future work should: (1) profile actual Qwen-32B and other vision-language model latencies on H100 hardware under production constraints, (2) categorize edits by structural complexity to identify whether large-error neurons require proportionally more inference passes, (3) measure how proofreading corrections compound (i.e., whether fixing one error reduces downstream corrections), and (4) validate projections through pilot deployments on held-out test datasets.

Appendix: Supplementary Analysis

GPU Computational Cost Analysis: Operational Effort and Computational Burden Across Species

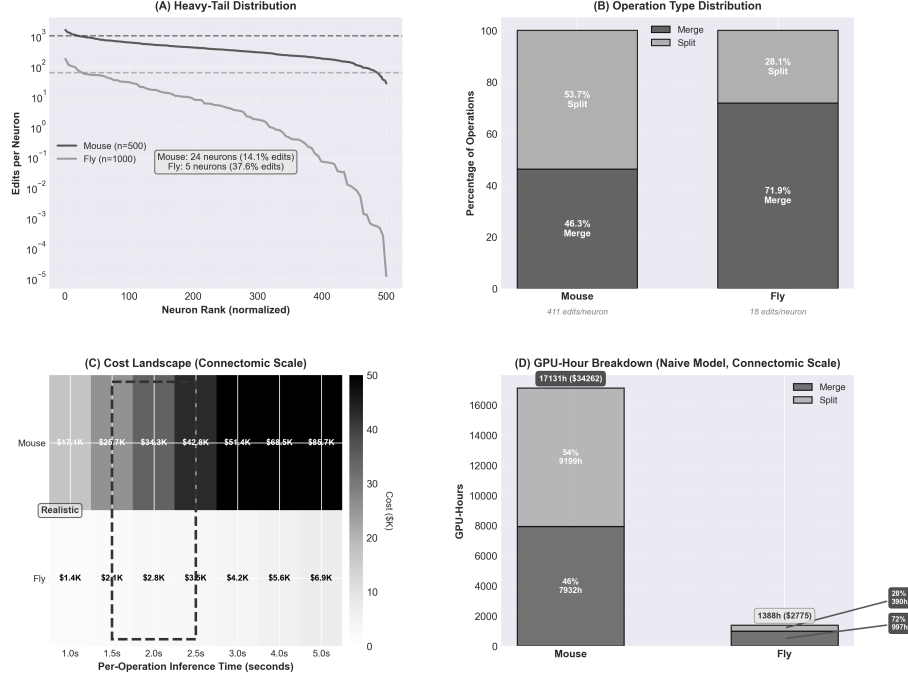
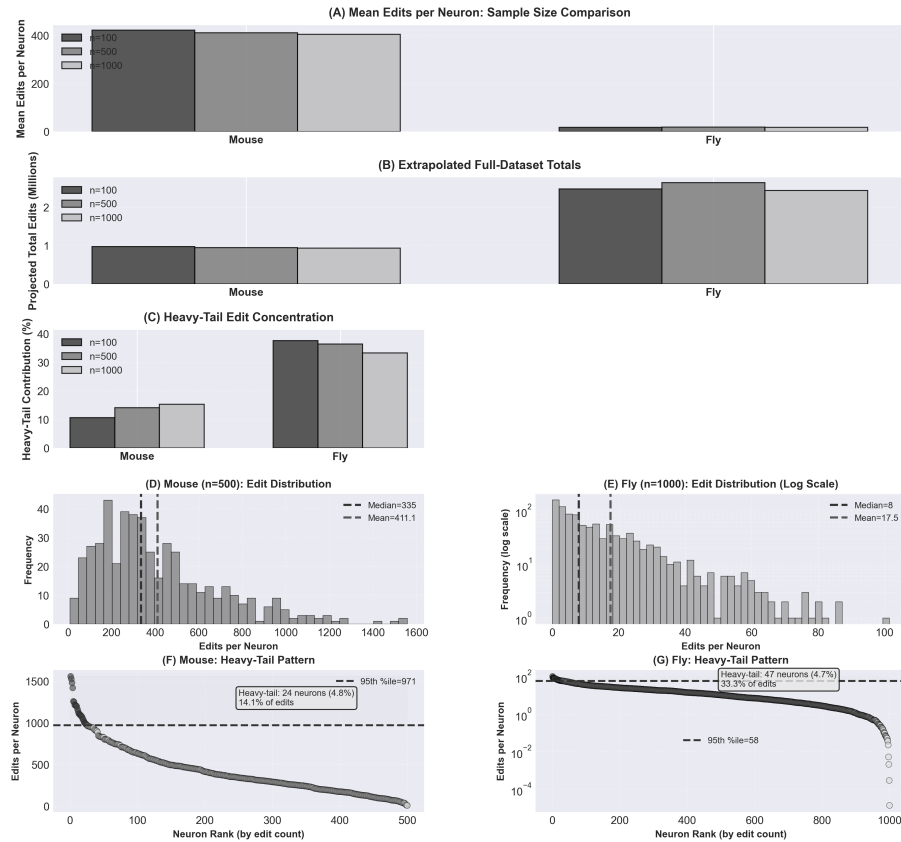


Figure 1: GPU Computational Cost Analysis: Operational Effort and Computational Burden Across Species. Main 4-panel figure combining heavy-tail distribution patterns (top) with computational cost analysis at connectomic scale (bottom). **(A) Heavy-Tail Edit Distribution.** Rank-ordered cumulative edit counts for mouse ($n = 500$, dark gray) and fly ($n = 1000$, light gray) on log scale, normalized to same x-axis for visual comparison. Dashed horizontal lines indicate 95th percentile thresholds; inset annotation box reports heavy-tail concentration statistics. Mouse heavy-tail neurons (24 neurons, 4.8% of sample, > 971 edits) contribute only 14.1% of total edits, indicating relatively uniform workload. Fly heavy-tail neurons (47 neurons, 4.7% of sample, > 58 edits) concentrate 33.3% of total edits, revealing that fly proofreading workload is dominated by a small number of complex outliers—a $23.4\times$ difference in per-neuron baseline despite $2.4\times$ higher per-neuron concentration. **(B) Operation Type Distribution.** Stacked bar chart showing merge vs. split operation percentages for mouse (46.3% merge, 53.7% split) and fly (71.9% merge, 28.1% split) across full proofread populations (2,314 and 139,255 neurons, respectively). Mouse’s balanced merge/split ratio indicates comparable over- and under-segmentation errors in the initial volume; fly’s merge-dominance reveals systematic oversegmentation in FAFB volume, requiring fragment consolidation as the dominant proofreading task. Mean edits per neuron noted below each species. **(C) Cost Landscape Heatmap.** 2D cost matrix (grayscale, vmax clipped to 50K to preserve visibility across range) showing GPU cost dependency on per-operation inference time (1.0–5.0 seconds) at connectomic volume scale (75,000 mouse neurons in $\tilde{P}mm^3$; 140,000 fly neurons in full brain). Dashed rectangle highlights the realistic 1.5–2.5 second range derived from task-stratified complexity analysis. White text on dark (mouse) and black text on light (fly) backgrounds ensure readability across the grayscale gradient. Cost estimates shown in \$K (thousands). **(D) GPU-Hour Breakdown (Naive Model, Connectomic Scale).** Stacked bar chart quantifying merge vs. split GPU-hours required for full-volume proofreading using naive uniform 2.0 s per-operation model. Mouse: 9,199 merge hours (54%) + 7,937 split hours (46%)

Supplementary Figure: Sample Size Validation & Distribution Analysis



Panels A–C: Sample Size Validation and Robustness. (A) Mean Edits per Neuron. Per-neuron edit statistics across increasing sample sizes ($n=100$, $n=500$, $n=1000$) for both species, using grayscale bars. Demonstrates robust sampling stability: mouse mean stabilizes around 411 edits/neuron with $< 2\%$ variation; fly stabilizes around 17.5 edits/neuron with $< 8\%$ variation, well within acceptable bounds for statistical inference. Validates that samples are representative of full proofread populations. **(B) Extrapolated Full-Dataset Totals.** Projected total edits when extrapolating samples to their respective full proofread populations (2,314 mouse neurons in 1mm³ MICrONS volume; 139,255 neurons in complete FlyWire brain). Shows convergence and stability of extrapolation with increasing sample size, demonstrating consistency of the linear extrapolation methodology used in main figure calculations. **(C) Heavy-Tail Edit Concentration.** Percentage of total edits contributed by heavy-tail neurons (defined as > 95 th percentile threshold) across different sample sizes. Validates that heavy-tail concentration patterns are consistent and robust across sample sizes, with mouse $\sim 14\%$ and fly $\sim 33\%$ contributions regardless of n .

Panels D–G: Edit Distribution Patterns and Heavy-Tail Structure Analysis. (D) Mouse ($n=500$) Edit Distribution Histogram.

Grayscale histogram of edit counts across 500 sampled mouse neurons on linear scale, showing pronounced right-skewed distribution. Median (335 edits, dark dashed line) substantially below mean (411.1 edits, gray dashed line), indicating strong heavy-tail effect. Concentration of neurons at low edit counts (0–200 range) emphasized by linear y-axis scaling. **(E) Fly ($n=1000$) Edit Distribution Histogram (Log Scale).** Grayscale histogram of edit counts across 1,000 sampled fly neurons using log-scale y-axis to visualize the full