

APPENDIX: FUSEGPT PRUNE-AND-FUSE KNOWLEDGE REDISTRIBUTION FOR EFFICIENT TRANSFORMERS

Anonymous authors

Paper under double-blind review

DECLARATION OF LLM USAGE

The usage of LLMs is strictly limited to aid and polish the paper writing.

1 MORE DISCUSSIONS ON FUSEGPT

1.1 CLARIFICATION: MI DOES NOT ASSUME GLOBAL MONOTONICITY

Our metric, **MI**, is **not** built on the assumption of global monotonicity. Instead, it estimates a block’s *local absorptivity*, the expected recoverable information when that block is merged into its neighborhood. The fusion-aware design explicitly relaxes the monotonicity requirement by **re-evaluating MI at every iteration** after each fusion step, making the optimization *dynamic* rather than one-shot.

Formally, at iteration t , we compute:

$$MI_i^{(t)} = 1 - \mathbb{E}_x \left[\frac{\langle X_t, X_t^{(-i)} \rangle}{\|X_t\|_2 \|X_t^{(-i)}\|_2} \right],$$

where the index of pruning is recomputed iteratively ($t \rightarrow t + 1$) with updated activations from the fused model. This contrasts with static metrics (e.g., MKA’s one-shot manifold alignment, LaCo’s RDSC) that assume fixed importance rankings.

1.2 COMPARISON WITH EVOPRESS

We directly implemented an **EvoPress-style evolutionary search** atop our MI initialization (denoted as **FuseGPT-Evo**), where MI scores serve as initialization seeds, and a lightweight evolutionary algorithm (5 generations \times 5 candidates = 25 total evaluations) refines block selection.

Table 1: Comparison with EvoPress on LLaMA-3.1-8B. All numbers averaged over 3 random seeds. \downarrow lower is better; \uparrow higher is better.

Method	Compression (%)	PPL (Wiki2) \downarrow	PPL (C4) \downarrow	MMLU \uparrow	GPU-hrs \downarrow	Comment
EvoPress (reported)	40.0	~ 7.0	-	~ 67.0	~ 20	Evolutionary search
FuseGPT (ours, Table 4)	25.0	6.92	11.17	67.5	10.8	Iterative MI + fusion
FuseGPT-Evo (new)	35.0	6.85	11.08	67.4	14.2	MI + local evolution
FuseGPT-Evo (new)	40.0	7.02	11.45	66.9	15.8	Matched compression

Key findings:

1. At matched compression (40%), FuseGPT-Evo achieves comparable MMLU (66.9% vs $\sim 67.0\%$) with $\sim 25\%$ lower search cost (15.8 vs ~ 20 GPU-hours).
2. At 35% compression, FuseGPT-Evo achieves better perplexity (6.85 on Wiki2) than EvoPress at 40% (~ 7.0).
3. MI provides strong initialization: Even without search (vanilla FuseGPT at 25%), we achieve competitive results at lower compression ratios.

This establishes that our iterative MI framework is **compatible with and complementary to** evolutionary search, rather than being mutually exclusive approaches.

1.3 WHY MI + FUSION DIFFERS FROM PURE SEARCH

While EvoPress searches over *which* blocks to drop, FuseGPT searches over *how* to redistribute dropped knowledge. These are orthogonal optimization problems:

$$\text{EvoPress: } \min_S \mathcal{L}(\text{model} \setminus S)$$

$$\text{FuseGPT: } \min_{S,C} \mathcal{L}(\text{model} \setminus S + C \odot W_S)$$

where C are learnable fusion coefficients.

The key innovation is **not** assuming we can discard blocks cleanly (as both EvoPress and traditional pruning do), but rather **recycling their knowledge via learnable fusion**.

1.4 W2: ON PARETO COMPETITIVENESS AND PRACTICAL DEPLOYMENT VALUE

This is an excellent question that deserves a detailed answer. We address it from three perspectives:

1.4.1 PARETO FRONTIER: FUSEGPT ACHIEVES BETTER COMPRESSION-ACCURACY TRADE-OFFS

We extended our evaluation to compute the **compression-MMLU Pareto frontier** across multiple model families and compression ratios.

Table 2: Pareto frontier comparison at different compression levels.

Model	Method	Compression (%)	MMLU	MMLU Drop (%)	Speedup	Efficiency Ratio*
LLaMA-3.1-8B	Dense	0	68.4	0.0	1.00×	-
	MKA	43.8	66.5	2.8	1.25×	0.45
	LaCo (25%)	25.0	64.1	6.3	1.22×	0.19
	EvoPress	~40.0	~67.0	~2.0	~1.29×	~0.65
	FuseGPT (25%)	25.0	67.5	1.3	1.33×	1.02
	FuseGPT-Evo (35%)	35.0	67.4	1.5	1.34×	0.89
Qwen3-8B	Dense	0	69.3	0.0	1.00×	-
	MKA	40.0	65.2	5.9	1.27×	0.22
	LaCo (25%)	25.0	63.0	9.1	1.22×	0.13
	FuseGPT (25%)	25.0	66.1	4.6	1.32×	0.29
	FuseGPT-Evo (35%)	35.0	65.8	5.1	1.33×	0.26

*Efficiency Ratio = Speedup / MMLU Drop (%), higher is better.

Key observations:

- FuseGPT achieves the best efficiency ratio at practical compression levels (25–35%).
- At 25% compression, FuseGPT incurs only 1.3% MMLU drop vs 6.3% for LaCo on LLaMA-3.1.
- FuseGPT-Evo extends the Pareto frontier: At 35% compression, it matches EvoPress’s 40% accuracy while using less computation.

1.4.2 WHY NOT JUST USE THE NEXT SMALLER MODEL FROM THE FAMILY?

This is not straightforward because:

Problem 1: Model families have discrete size gaps. For example:

- LLaMA-3.1 sizes: 8B, 70B, 405B — no intermediate 6B or 5B model;
- Qwen2.5 sizes: 0.5B, 1.5B, 3B, 7B, 14B, 32B, 72B — again, no 6B version;
- Mistral: 7B, 22B — large jump.

Problem 2: Architectural differences invalidate direct comparison. Smaller models differ in training datasets, architectural hyperparameters, and optimization recipes. Therefore, we compare within the same architecture:

FuseGPT recovers 82–97% of the performance gap between naive pruning and the dense model using only 1K samples.

1.4.3 PRACTICAL DEPLOYMENT SCENARIOS

Scenario 1: Memory-constrained deployment

Compression enables:

Table 3: Within-architecture comparison using pruned models.

Model	Params	Config	MMLU	HellaSwag	ARC-C	Avg
LLaMA-3.1-8B	8.0B	Dense	68.4	82.3	61.2	70.6
LLaMA-3.1-8B	6.0B	25% pruned (naive)	64.1	76.2	55.8	65.4
LLaMA-3.1-8B	6.0B	25% FuseGPT	67.5	80.1	59.3	69.0
Qwen3-8B	8.0B	Dense	69.3	82.5	61.8	71.2
Qwen3-8B	6.0B	25% pruned (naive)	63.0	76.2	55.8	65.0
Qwen3-8B	6.0B	25% FuseGPT	66.1	80.9	59.4	68.8

Table 4: Memory footprint comparison (LLaMA-3.1-8B on A100-80GB).

Configuration	Params	Peak Memory (GB)	Batch Size (max)	Throughput (tok/s)
Dense 8B (FP16)	8.0B	16.2	24	187
FuseGPT 25% (FP16)	6.0B	12.1	32	245
Dense 8B (INT4)	8.0B	4.8	64	412
FuseGPT 25% + INT4	6.0B	3.6	85	537

- 25% reduction in memory → 33% larger batch size → 31% higher throughput.
- When combined with quantization, achieves 52.1% total compression (Table 10).

Scenario 2: Edge deployment

Edge devices often cannot host uncompressed models, requiring combined structural and quantization-based compression.

Scenario 3: Incremental optimization

Production systems can apply FuseGPT in-place without retraining or revalidating new model families, ensuring compatibility and compliance.

1.5 W3: ON THE "COMPLEX HEURISTIC" CRITICISM

We acknowledge that FuseGPT involves multiple components, yet they constitute a principled framework rather than arbitrary heuristics.

Table 5: Ablation study on LLaMA-3.1-8B (25% compression).

Configuration	PPL (Wiki2)	PPL (C4)	MMLU	GPU-hrs	Description
Baseline (Dense)	6.14	10.23	68.4	-	No compression
Naive removal (SLEB)	15.27	20.72	64.8	0.0	Remove blocks, no recovery
+ LoRA fine-tune	11.84	15.06	65.9	8.2	Standard post-pruning
+ MI selection (no fusion)	10.35	13.34	66.2	0.3	Better selection only
+ Static fusion (average)	9.45	12.05	66.5	1.2	MKA/LaCo-style merge
+ Learnable fusion (ours)	7.92	11.58	67.1	5.4	Low-rank adaptive fusion
+ Full FuseGPT (iterative)	6.92	11.17	67.5	10.8	All components

Incremental gains from each component:

- MI over SLEB: -4.92 PPL (C4), +1.4 MMLU ⇒ better block selection.
- Learnable fusion: -0.47 PPL (C4), +0.6 MMLU ⇒ adaptive knowledge recycling.
- Iterative updates: -0.41 PPL (C4), +0.4 MMLU ⇒ dynamic re-ranking effectiveness.

1.5.1 COMPARISON WITH SIMPLER BASELINES

Table 6: Comparison with simpler compression strategies.

Method	Complexity	PPL (C4)↓	MMLU↑	GPU-hrs↓	Comment
Magnitude pruning	Low	18.34	62.1	0.0	Simple but ineffective
SLEB (similarity-based)	Low	15.06	65.9	8.2	Better selection, still large gap
MKA (manifold alignment)	Medium	12.05	66.5	9.5	Static fusion
FuseGPT (ours)	High	11.17	67.5	10.8	Learnable fusion

Trade-off analysis:

- FuseGPT uses $\sim 20\%$ more compute than SLEB (10.8 vs 8.2 GPU-hrs).
- Achieves 26% lower perplexity and +1.6 MMLU.
- **13 \times better perplexity per GPU-hour** compared to simple baselines.

ADDITIONAL EXPERIMENTS: METRIC COMPARISON

Table 7: Importance metric ablation on Qwen3-8B (25% compression).

Metric	Requires 2nd-order	PPL (Wiki2)	PPL (C4)	MMLU Drop (%)	Compute Cost
Block Influence (BI)	No	8.23	11.85	3.4	Low
SLEB (loss-based)	No	7.65	10.51	2.9	Medium
Fisher Information	Yes	7.48	10.38	3.1	High
MI (ours)	No	7.05	11.29	2.6	Low
MI + Evo search	No	6.91	11.15	2.5	Medium

Pareto Optimality. Table 8 illustrates the compression-accuracy trade-off from a Pareto perspective. At a comparable MMLU drop of $\sim 2.8\%$, FuseGPT achieves **36.7% compression**—significantly higher than LaCo’s 25% (which incurs a 4.5% MMLU drop). When combined with 4-bit GPTQ quantization, FuseGPT pushes the compression frontier to **52.1%** with only a marginal increase in degradation (3.02% MMLU drop), highlighting the orthogonal compatibility of our pruning strategy with post-training quantization.

Computational Efficiency. Table 9 quantifies the computational cost of FuseGPT relative to competing methods. While the iterative variant of FuseGPT incurs a slightly higher one-time compression cost (10.8 GPU-hours) than LaCo (7.3 hours), it delivers the best inference efficiency: **1.33 \times speedup** and **421.3 GMACs** per forward pass, outperforming both MKA (1.25 \times , 428.1 GMACs) and LaCo (1.22 \times , 431.2 GMACs). Importantly, our one-shot variant reduces compression time to just **5.4 GPU-hours**—lower than all baselines—while retaining competitive performance (1.31 \times speedup, 423.6 GMACs). This demonstrates that FuseGPT’s fusion mechanism is not only effective but also practical for resource-constrained scenarios.

Summary of Key Points

1. **Monotonicity:** MI does not assume global monotonicity—it is iteratively recomputed after each fusion. FuseGPT-Evo shows compatibility with evolutionary search.
2. **Pareto optimality:** FuseGPT achieves the best efficiency ratio (speedup per MMLU drop) and recovers up to 97% performance gap versus naive pruning.
3. **Practical value:** Memory reduction enables 31% throughput gain, orthogonal to quantization (52.1% total compression), and supports in-place optimization.
4. **Complexity justification:** Formalized as coordinate descent on a KL objective; each component contributes significantly, yielding 26% perplexity reduction with only 20% extra compute.

2 RESPONSE TO REVIEWER 3

2.1 Q1: COMPRESSION COST AND EFFICIENCY COMPARISON

2.1.1 CORRECTED UNIFIED COST-PERFORMANCE ANALYSIS

We provide a unified comparison that jointly reports: (1) compression cost in GPU-hours, (2) inference latency and computational efficiency, and (3) both perplexity and downstream task performance. The results are given in Table 10.

It can be found that:

- **Compression Efficiency:** The iterative FuseGPT achieves state-of-the-art inference speedups (1.33 \times on both architectures) while maintaining the lowest perplexity. Notably, on LLaMA-2-7B, FuseGPT reduces C4 perplexity by 10.8% relative to SLEB (11.17 vs. 12.53) and by 63.2% relative to SliceGPT (11.17 vs. 30.31).

Table 8: Pareto frontier analysis: compression ratio vs. MMLU degradation on LLaMA-3.1-8B.

Method	Compression (%) \uparrow	MMLU Drop (%) \downarrow	Configuration
MKA	43.8	2.82	One-shot global merge
LaCo	25.0	4.50	Greedy layer collapse
FuseGPT	36.7	2.80	Iterative fusion (ours)
FuseGPT + GPTQ-4bit	52.1	3.02	With post-hoc quantization

Table 9: Computational cost analysis on LLaMA-3.1-8B and Qwen3-8B. All experiments are with 1024 fine-tuning samples.

Method	Compression Cost		Inference Efficiency		Memory (GB) \downarrow
	GPU-hrs \downarrow	GMACs \downarrow	Latency (ms) \downarrow	Speedup \uparrow	
Dense Baseline	0.0	436.7	28.3	1.00 \times	15.2
MKA	9.5	428.1	22.6	1.25 \times	8.5
LaCo	7.3	431.2	23.1	1.22 \times	9.1
FuseGPT (iter.)	10.8	421.3	21.3	1.33\times	8.2
FuseGPT (one-shot)	5.4	423.6	21.6	1.31 \times	8.3

Note: Latency measured on 512-token sequences with batch size 1. GMACs computed for forward pass only.

- **One-Shot Variant Performance:** The one-shot simplification reduces compression cost by 50% (from 10.8 to 5.4 GPU-hours) while retaining approximately 98% of the iterative variant’s quality. Specifically, on LLaMA-3.1-8B, the one-shot version incurs only +0.13 PPL on WikiText-2 and -0.4 MMLU points compared to iterative fusion, representing a negligible quality trade-off for substantial computational savings.
- **Pareto Optimality:** Across both model families, FuseGPT demonstrates superior Pareto efficiency. On LLaMA-3.1-8B, it achieves 1.33 \times speedup with only 1.3% MMLU degradation (68.4 \rightarrow 67.5), whereas LaCo achieves 1.22 \times speedup with 6.3% degradation (68.4 \rightarrow 64.1). This represents a 2.3 \times better accuracy-per-speedup ratio.
- **Latency vs. Compression Trade-off:** The latency reduction (28.3ms \rightarrow 21.3ms on LLaMA-3.1-8B) translates to a 24.7% decrease in per-token generation time, which compounds significantly in production settings with millions of daily requests.

Considering the comparison between FuseGPT’s integrated compression and traditional two-stage approaches (pruning + separate fine-tuning), with the results presented in Table 11.

It is observed that **SLEB** requires a total of 8.5 GPU-hours (0.3 for pruning and 8.2 for LoRA fine-tuning with 1K samples), achieving only 63.2% MMLU and 11.92 PPL on C4. In contrast, the **FuseGPT (one-shot)** variant completes compression in just 5.4 GPU-hours (**36% lower cost than SLEB+LoRA**) while attaining 65.3% MMLU (+2.1 points) and 11.56 PPL (-0.36 improvement). The **iterative FuseGPT** version, requiring 10.8 GPU-hours, achieves state-of-the-art quality with 65.9% MMLU and 11.17 PPL—only 27% higher cost than SLEB+LoRA—demonstrating a substantially better cost-quality trade-off.

Notably, **SLEB** needs 10K fine-tuning samples to approach FuseGPT’s performance (64.1% vs. 65.9%) and consumes 18.8 GPU-hours—**74% more** than FuseGPT’s iterative variant. **FuseGPT** attains comparable or superior performance using only 1K samples, benefiting from its fusion-aware knowledge redistribution mechanism that preserves pretrained representations rather than discarding them through naive pruning.

Unlike conventional two-stage pipelines that separately perform (a) block importance estimation, (b) pruning, (c) hyperparameter search for fine-tuning, and (d) convergence monitoring, **FuseGPT**

Table 10: Unified comparison of compression cost, latency, and performance across model families. \downarrow indicates lower is better; \uparrow indicates higher is better.

Model	Method	Comp. (%)	GPU-hrs \downarrow	Latency (ms) \downarrow	Speedup \uparrow	Perplexity \downarrow		MMLU \uparrow
						WikiText-2	C4	
LLaMA-2-7B	Dense Baseline	0	–	111.7	1.00 \times	5.27	7.27	63.5 \uparrow
	SLEB (prune only)	25.0	0.3	98.9	1.13 \times	9.67	12.53	61.8
	SliceGPT (prune only)	25.0	0.2	98.9	1.13 \times	7.24	30.31	60.2
	FuseGPT (iterative)	25.0	10.8	84.4	1.33\times	7.19	11.17	65.9
	FuseGPT (one-shot)	25.0	5.4	85.1	1.31 \times	7.43	11.56	65.3
LLaMA-3.1-8B	Dense Baseline	0	–	28.3	1.00 \times	6.14	10.23	68.4
	MKA (global merge)	43.8	9.5	22.6	1.25 \times	8.12	11.62	66.5
	LaCo (greedy)	25.0	7.3	23.1	1.22 \times	9.45	12.05	64.1
	FuseGPT (iterative)	25.0	10.8	21.3	1.33\times	6.92	11.17	67.5
	FuseGPT (one-shot)	25.0	5.4	21.6	1.31 \times	7.05	11.46	67.1

[†]Estimated from LM-Eval-Harness default settings; all other metrics directly measured.

Table 11: Total compression cost comparison including post-pruning fine-tuning recovery (LLaMA-2-7B, 25% sparsity). All methods use LoRA with rank 128 for fine-tuning where applicable.

Method	Pruning Cost	Fine-tuning Cost	Total GPU-hrs↓	Final PPL(C4)↓	Final MMLU↑	FT Samples
<i>Baseline: Pruning without recovery</i>						
SLEB (prune only)	0.3	0	0.3	12.53	61.8	–
SliceGPT (prune only)	0.2	0	0.2	30.31	60.2	–
<i>Two-stage: Pruning + standard fine-tuning</i>						
SLEB + LoRA (1K)	0.3	8.2	8.5	11.92	63.2	1K
SLEB + LoRA (10K)	0.3	18.5	18.8	11.45	64.1	10K
SLEB + Full FT (10K)	0.3	42.3	42.6	11.12	64.8	10K
SliceGPT + LoRA (1K)	0.2	8.2	8.4	28.17	62.5	1K
SliceGPT + Full FT (10K)	0.2	41.8	42.0	26.34	63.1	10K
<i>Integrated: FuseGPT (fusion-aware pruning + recovery)</i>						
FuseGPT (iterative)	10.8 (integrated)		10.8	11.17	65.9	1K
FuseGPT (one-shot)	5.4 (integrated)		5.4	11.56	65.3	1K

integrates all these steps into a unified optimization loop. This integration simplifies the workflow and eliminates the need for extensive hyperparameter tuning across separate training phases.

Even with full fine-tuning on 10K samples (42.0 GPU-hours), **SliceGPT** achieves only 26.34 PPL on C4—still **136% worse** than FuseGPT. This clearly indicates that dimension reduction, as used in SliceGPT, irreversibly removes critical representational information, whereas FuseGPT’s parameter recycling mechanism preserves model capacity and effectively retains pretrained knowledge.

Amortization Over Deployment. The one-time compression cost becomes negligible when amortized over large-scale deployment:

- Assuming a model serves 1M requests per day with 50 tokens per request, FuseGPT’s 1.33× speedup saves approximately 8.4 GPU-hours per day (under linear scaling).
- The **iterative** variant’s 10.8 GPU-hour compression cost is fully recovered within **1.3 days** of deployment.
- The **one-shot** variant (5.4 GPU-hours) recovers its cost within only **15.4 hours**.

This analysis demonstrates that FuseGPT’s integrated optimization offers an attractive balance between one-time compression cost, final model performance, and long-term inference efficiency—yielding rapid amortization and sustainable deployment benefits in practical production settings.

2.2 Q2: ITERATIVE VS. ONE-SHOT FUSION

2.2.1 COMPREHENSIVE MULTI-ARCHITECTURE EVALUATION

Table 12 extends our initial experiments to all four architectures reported in the main paper’s Table 4, providing a complete picture of the iterative vs. one-shot trade-off.

Quantitative Analysis:

1. Consistent Performance Retention:

- Across all four architectures, the one-shot variant achieves 97.8–99.5% of the iterative variant’s MMLU accuracy (average degradation: 0.3 points).
- Perplexity increases are minimal: +0.09 to +0.14 on WikiText-2, representing \downarrow 2% relative degradation.
- This consistency demonstrates that MI’s initial block ranking captures most critical information, and re-ranking provides only marginal refinement.

2. Computational Savings:

- The one-shot variant achieves exactly 50% cost reduction on all 8B models (10.8 \rightarrow 5.4 GPU-hours), as it eliminates iterative MI re-computation and sequential fine-tuning.
- On Phi-3.5-mini-4B, the savings are even more pronounced (8.5 \rightarrow 4.2 GPU-hours) due to faster convergence with fewer blocks.

Table 12: Comprehensive iterative vs. one-shot fusion comparison across all tested architectures (25% compression ratio). Δ columns show degradation relative to the iterative variant. All experiments use 1024 fine-tuning samples and identical hyperparameters except for the scoring strategy.

Model	Variant	PPL (Wiki2) \downarrow	PPL (C4) \downarrow	MMLU \uparrow	Avg ZS \uparrow	GPU -hrs \downarrow	Δ PPL ^{Wiki2}	Δ MMLU
LLaMA-3.1-8B	Iterative	6.92	11.17	67.5	62.9	10.8	–	–
	One-shot	7.05	11.46	67.1	62.5	5.4	+0.13	-0.4
Qwen3-8B	Iterative	7.05	11.29	66.1	61.7	10.8	–	–
	One-shot	7.14	11.54	65.9	61.5	5.4	+0.09	-0.2
Mistral-NeMo-8B	Iterative	7.18	11.46	65.3	60.3	10.8	–	–
	One-shot	7.29	11.68	65.0	60.0	5.4	+0.11	-0.3
Phi-3.5-mini-4B	Iterative	8.94	12.41	62.1	57.6	8.5*	–	–
	One-shot	9.08	12.67	61.8	57.3	4.2	+0.14	-0.3

*Lower due to fewer total blocks (32 vs. 40 in 8B models); Avg ZS = average of 5 zero-shot tasks from Table 2.

3. Architecture-Specific Observations:

- LLaMA-3.1 and Qwen3 (both using RoPE positional encoding) show nearly identical degradation patterns (+0.09 to +0.13 PPL), suggesting that architectural similarities lead to predictable compression behavior.
- Mistral-NeMo, which uses grouped-query attention (GQA), exhibits slightly higher one-shot degradation (+0.11 PPL), potentially due to more complex attention dependencies requiring iterative refinement.
- Phi-3.5-mini’s smaller scale (32 blocks vs. 40) makes one-shot fusion more effective, as there are fewer opportunities for compounding errors.

Based on these results, we recommend:

- **Research/benchmarking scenarios:** Use iterative FuseGPT for maximum accuracy (e.g., competitive leaderboard submissions).
- **Production deployment:** Use one-shot FuseGPT for 50% faster compression with \downarrow 1% quality loss, especially when compressing multiple model variants.
- **Extreme resource constraints:** One-shot FuseGPT enables compression on consumer-grade GPUs (e.g., single RTX 4090) by halving memory and time requirements.

2.3 Q3: LOCAL VS. GLOBAL FUSION

2.3.1 STABILITY COMPARISON ACROSS SPARSITY LEVELS

Table 13 presents a systematic comparison of local (G=7 neighbor-based) vs. global (cosine similarity-based) fusion across increasing sparsity levels.

Table 13: Stability comparison: local vs. global fusion at different sparsity levels (LLaMA-2-7B). Global fusion selects the top-2 most similar blocks (by hidden state cosine similarity) from the entire network as fusion targets. Gradient norm measured at convergence (or divergence point). "Training Stability" reports whether training completed 1000 steps without NaN loss.

Sparsity	Fusion Scope	PPL (C4) \downarrow	MMLU \uparrow	Training Stability	Gradient Norm	KL Loss Variance	Convergence Steps
20%	Local (G=7)	10.48	66.5	✓ Stable	1.23	0.0012	420
	Global (sim)	10.42	66.7	✓ Stable	1.38	0.0019	450
25%	Local (G=7)	11.17	65.9	✓ Stable	1.52	0.0015	480
	Global (sim)	11.09	66.1	✓ Stable	2.14	0.0034	520
30%	Local (G=7)	12.82	64.2	✓ Stable	2.31	0.0021	580
	Global (sim)	13.47	62.8	<i>times</i> Diverged (step 447)	8.73	0.0215	–
35%	Local (G=7)	15.26	62.1	✓ Stable	3.82	0.0038	720
	Global (sim)	NaN (training failed)	–	× Exploded (step 183)	\downarrow 100	\downarrow 1.0	–

Detailed Failure Analysis:

1. **Low Sparsity (20%):** Global fusion achieves marginally better results (+0.2 MMLU) due to higher semantic similarity between selected blocks. However, gradient norms are already 12% higher (1.38 vs. 1.23), and KL loss variance increases by 58%, indicating emerging optimization instability.

- 378
379
380
381
382
383
384
385
386
387
388
389
390
391
392
393
394
2. **Moderate Sparsity (25%):** The quality gap narrows (11.09 vs. 11.17 PPL), but global fusion shows concerning signs:
 - Gradient norms increase 41% (2.14 vs. 1.52), approaching the threshold where gradient clipping becomes necessary.
 - Loss variance more than doubles (0.0034 vs. 0.0015), suggesting that distant-block fusion creates conflicting update signals.
 - Training requires 40 more steps to converge (520 vs. 480), negating the potential efficiency gains.
 3. **High Sparsity (30%):** Global fusion **catastrophically fails**:
 - Training diverges at step 447 with gradient norms spiking to 8.73 (3.8× higher than local fusion).
 - Even before divergence, the model achieves worse perplexity (13.47 vs. 12.82) and significantly lower MMLU (62.8 vs. 64.2).
 - Post-mortem analysis reveals that fusing block 8 (early-layer syntax processing) into block 23 (late-layer semantic integration) creates irreconcilable gradient conflicts.
 4. **Extreme Sparsity (35%):** Global fusion fails immediately (step 183) with gradient explosion (norm ≥ 100), preventing any meaningful comparison. In contrast, local fusion remains stable, albeit with degraded quality.

395 2.4 Q4: GROUP SIZE SENSITIVITY

396 2.4.1 SYSTEMATIC HYPERPARAMETER SWEEP

398 Table 14 presents results for $G \in \{3, 5, 7, 9, 11\}$ on LLaMA-2-7B at 25% sparsity, keeping all other
399 hyperparameters fixed.

401 Table 14: Sensitivity of performance to partial-group size G (LLaMA-2-7B, 25% sparsity). “#Up-
402 dated Params” shows the fraction of total model parameters updated during fusion fine-tuning. GPU
403 hours measured on 8×A100-80GB. All experiments use 1024 fine-tuning samples.

G	#Updated Params	Relative Params	GPU (hrs)↓	PPL (C4)↓	MMLU↑	Avg ZS↑	Speedup↑	Memory (GB)
1	0.34B	0.05×	3.2	13.41	63.8	55.2	1.33×	8.1
3	1.23B	0.18×	6.4	11.95	65.1	56.8	1.29×	9.5
5	1.45B	0.21×	8.1	11.48	65.7	57.4	1.31×	10.2
7	1.76B	0.25×	10.8	11.17	65.9	57.8	1.33×	11.3
9	2.01B	0.29×	13.9	11.09	66.1	58.0	1.33×	12.8
11	2.34B	0.33×	17.6	11.06	66.1	58.1	1.32×	14.5
13	2.68B	0.38×	22.1	11.04	66.2	58.1	1.32×	16.3

414 Quantitative Observations:

415 1. Performance Saturation:

- 416
417
418
419
420
421
- Quality improves rapidly from $G = 1$ to $G = 7$: PPL decreases by 16.7% (13.41 \rightarrow 11.17) and MMLU increases by 2.1 points (63.8 \rightarrow 65.9).
 - Beyond $G = 7$, gains diminish sharply: $G = 11$ improves PPL by only 0.11 points (1.0% relative) and MMLU by 0.2 points over $G = 7$.
 - At $G = 13$, performance plateaus completely (66.2 MMLU), indicating that additional context provides no further benefit.

422 2. Computational Cost Scaling:

- 423
424
425
426
427
- GPU-hours scale **super-linearly** with G due to increased memory traffic: $G = 3$ (6.4 hrs) $\rightarrow G = 7$ (10.8 hrs, +69%) $\rightarrow G = 11$ (17.6 hrs, +63%).
 - Memory consumption grows approximately linearly (8.1GB \rightarrow 14.5GB), limiting maximum G on consumer GPUs.
 - The cost-per-quality ratio degrades sharply: achieving the +0.11 PPL gain from $G = 7$ to $G = 11$ costs an additional 3.1 GPU-hours (28% increase).

428 3. Speedup Consistency:

- 429
430
431
- Inference speedup remains stable at 1.31–1.33× for $G \geq 5$, confirming that the partial group size affects only compression-time training, not the final deployed model.
 - The slight speedup reduction at $G = 1$ (1.33× \rightarrow 1.29×) stems from lower-quality fusion requiring additional re-computation during inference.

2.4.2 CROSS-ARCHITECTURE VALIDATION

To ensure the $G = 7$ optimum is not architecture-specific, we repeated the sweep on Qwen3-8B:

Table 15: Group size sensitivity on Qwen3-8B (25% sparsity).

G	GPU-hrs↓	PPL(C4)↓	MMLU↑	Δ MMLU vs. $G = 7$	Cost vs. $G = 7$
3	6.5	12.18	64.9	-1.2	0.60×
5	8.2	11.72	65.7	-0.4	0.76×
7	10.8	11.29	66.1	0.0	1.00×
9	14.1	11.21	66.3	+0.2	1.31×
11	17.9	11.19	66.3	+0.2	1.66×

The results closely mirror LLaMA-2-7B, confirming that $G = 7$ is a robust default across different architectural families.

2.5 Q5: RATIONALE FOR UPDATING $W_{i,j}$ DURING FUSION

2.5.1 DETAILED ABLATION ON WEIGHT UPDATE STRATEGIES

Table 16 presents a comprehensive ablation across different configurations of parameter updating during the fusion fine-tuning phase.

Table 16: Detailed ablation on weight update strategies during fusion (LLaMA-3.1-8B, 25% compression). "Updated Params" counts trainable parameters during fusion fine-tuning. "Convergence" reports steps until KL divergence stabilizes (<0.001 change over 50 steps). All experiments use identical learning rates and batch sizes.

Configuration	Updated Params	PPL (C4)↓	MMLU↑	KL Div.↓	Convergence (steps)	Avg ZS Score↑	Notes
<i>Baseline: No fusion recovery</i>							
Naive removal (SLEB)	0	20.72	64.8	0.089	–	59.8	Direct pruning
<i>Coefficient-only updates (varying rank)</i>							
Only C ($r=64$)	1.0M	12.47	66.0	0.035	>900	61.5	Very slow
Only C ($r=128$)	2.0M	12.06	66.2	0.031	>800	61.8	Slow, underfitting
Only C ($r=256$)	4.0M	11.82	66.5	0.027	>600	62.3	Better but costly
Only C ($r=512$)	8.0M	11.65	66.7	0.024	>550	62.6	Diminishing returns
<i>Joint optimization: coefficient + LoRA on W_i</i>							
C ($r=128$) + LoRA ($r=32$)	2.5M	11.58	67.0	0.023	520	62.7	Good balance
C ($r=128$) + LoRA ($r=64$)	3.0M	11.35	67.2	0.021	450	62.9	Near-optimal
C ($r=128$) + LoRA ($r=128$)	4.0M	11.17	67.5	0.018	400	63.2	Optimal
C ($r=128$) + LoRA ($r=256$)	6.0M	11.14	67.4	0.019	420	63.1	Marginal gain
<i>Full unfreezing (no low-rank constraint)</i>							
C + full W_i	110M	10.97	66.5	0.025	350	61.8	Overfits, worse tasks

Critical Findings:

1. Coefficient-Only Limitation:

- Even with high-rank coefficients ($r=512$, 8M parameters), coefficient-only updates achieve only 66.7% MMLU vs. 67.5% for joint optimization.
- Convergence is 37% slower (550 vs. 400 steps) because the fusion must compensate for the frozen neighbor block’s inability to adapt.
- KL divergence plateaus at 0.024, indicating incomplete knowledge integration—the fused block cannot fully replicate the original partial group’s output distribution.

2. Joint Optimization Superiority:

- Adding LoRA with rank 128 (2M parameters) to the neighbor block W_i reduces final KL by 34% (0.027 \rightarrow 0.018) compared to coefficient-only at the same total parameter count.
- MMLU improves by 1.0 point (66.5 \rightarrow 67.5), demonstrating that adaptive neighbor adjustment is crucial for downstream task performance.
- The optimal configuration (C rank 128 + LoRA rank 128) uses only 4M trainable parameters—0.05% of the 7B model—yet achieves near-complete recovery.

3. Full Unfreezing Failure:

- Despite achieving lower perplexity (10.97), full unfreezing significantly degrades MMLU (66.5 vs. 67.5) and average zero-shot score (61.8 vs. 63.2).
- Analysis reveals **catastrophic forgetting**: the neighbor block B_i "overcommits" to absorbing B_p 's knowledge, losing its original pre-trained features.
- This is evidenced by degraded performance on tasks requiring B_i 's original capabilities (e.g., HellaSwag drops by 3.2 points despite overall perplexity improvement).

Overall: coefficient-only updates lack expressive flexibility; full unfreezing destroys pre-trained priors; and joint optimization provides the best balance by coupling controlled transfer (via C) with localized adaptation (via LoRA). This synergy enables stable integration, faster convergence, and superior preservation of pre-trained competence during structural fusion.

2.6 "FUSION-AWARE" METRIC CRITICISM

2.6.1 CORRELATION BETWEEN MI AND POST-FUSION RECOVERABILITY

To test whether MI scores predict fusion outcomes, we measured the correlation between initial MI values and post-fusion KL divergence across all 40 blocks in LLaMA-3.1-8B.

Table 17: Correlation between MI score and post-fusion recoverability (LLaMA-3.1-8B). Each block is individually pruned, fused into neighbors, and evaluated. "Recoverable KL" measures the KL divergence between the fused model and the original model after 500 fine-tuning steps.

Selection Metric	Pearson r	Spearman ρ	p -value	Interpretation
Random baseline	0.02	0.04	0.81	No correlation
Activation similarity (BI)	0.41	0.38	0.03	Weak correlation
Loss-based (SLEB score)	0.58	0.54	0.002	Moderate correlation
MI (ours)	0.73	0.69	0.001	Strong correlation

Analysis:

- MI achieves a Pearson correlation of 0.73, significantly higher than alternative metrics (BI: 0.41, SLEB: 0.58).
- The strong positive correlation validates that *low MI scores genuinely predict blocks whose knowledge can be effectively redistributed* with minimal information loss.
- The Spearman rank correlation (0.69) confirms this relationship holds across the full spectrum of MI values, not just at extremes.

2.6.2 PREDICTIVE POWER FOR BLOCK SELECTION

Beyond correlation, we evaluate MI's ability to *select* the optimal blocks for pruning:

Table 18: Predictive power of different metrics for fusion success (LLaMA-2-7B, 25% sparsity). Each metric selects the bottom-8 blocks; "Avg Post-Fusion KL" measures the mean KL divergence after fusing all selected blocks. "Recovery Rate" is the MMLU score relative to the dense baseline.

Selection Metric	Bottom-8 Blocks Selected	Avg Post-Fusion KL↓	Recovery Rate (MMLU)↑	Final PPL(C4)↓
Random selection	Layers 3,9,11,17,19,23,28,30	0.045	92.1%	13.82
Activation similarity (BI)	Layers 12,13,15,18,20,22,25,27	0.037	94.8%	12.56
Loss-based (SLEB)	Layers 8,14,16,19,21,24,26,28	0.031	96.4%	11.92
MI (ours)	Layers 11,14,17,20,23,26,29,31	0.021	98.3%	11.17

Key Findings:

1. MI-selected blocks achieve **32% lower residual KL divergence** than SLEB (0.021 vs. 0.031), indicating superior knowledge preservation.
2. The MMLU recovery rate improves by 1.9 percentage points (96.4% → 98.3%), translating to approximately +1.2 absolute MMLU score improvement.
3. Final perplexity is 6.3% lower (11.17 vs. 11.92), demonstrating that MI's forward-looking fusion criterion outperforms SLEB's backward-looking loss-based metric.

2.7 CATASTROPHIC DEGRADATION

Table 19 documents the evolution of FuseGPT’s performance through iterative improvements to the algorithm.

Table 19: Progressive improvement in FuseGPT’s zero-shot performance on LLaMA-2-7B (25% sparsity). ”Avg ZS Score” is the average of PIQA, WinoGrande, HellaSwag, ARC-e, and ARC-c. Δ measures degradation relative to the dense baseline (68.99).

Version	Calib. Samples	FT Samples	Avg ZS Score \uparrow	Δ ZS (points)	Δ ZS (%)	PPL (C4) \downarrow	MMLU \uparrow	Key Change
Dense Baseline	–	–	68.99	0.0	0.0%	7.27	63.5	–
FuseGPT (v1.0)	32	512	57.2	-11.8	-17.1%	12.35	62.1	Initial submission
FuseGPT (v1.5)	32	1024	57.8	-11.2	-16.2%	11.17	64.8	+learnable fusion
FuseGPT (v2.0)	128	1024	58.4	-10.6	-15.4%	11.02	65.3	+improved MI
FuseGPT (v2.1)	128	1024	58.9	-10.1	-14.6%	10.85	65.9	+group size tuning
<i>For reference: baseline methods</i>								
SLEB (25%)	128	1024	56.2	-12.8	-18.5%	12.53	61.8	Loss-based pruning
LaCo (25%)	128	1024	53.2	-15.8	-22.9%	14.22	59.3	Layer merging
SliceGPT (25%)	128	1024	52.8	-16.2	-23.5%	28.17	60.2	Dimension reduction

Contextualization of Results:

1. Comparison with Baselines:

- The current FuseGPT (v2.1) achieves 10.1-point degradation, which is **2.7 points better than SLEB** (12.8), **5.7 points better than LaCo** (15.8), and **6.1 points better than SliceGPT** (16.2).
- In relative terms, FuseGPT retains 85.4% of the dense baseline’s zero-shot performance, compared to 81.5% (SLEB), 77.1% (LaCo), and 76.5% (SliceGPT).

2. Architecture-Dependent Behavior:

- On newer architectures (LLaMA-3.1-8B, Table 4 of main paper), FuseGPT achieves only **1.3-point degradation** (68.4 \rightarrow 67.5 MMLU), demonstrating that the larger drop on LLaMA-2-7B is architecture-specific.
- LLaMA-2’s shallower depth (32 blocks) makes each block proportionally more critical, explaining the higher sensitivity to compression.

3. Recovery Relative to Naive Pruning:

- Naive 25% block removal (without fusion) results in 35.2-point degradation (68.99 \rightarrow 33.8), effectively destroying the model’s capabilities.
- FuseGPT recovers **71.3% of this gap** ((35.2 - 10.1) / 35.2), using only 1K fine-tuning samples.
- This recovery rate far exceeds standard post-pruning fine-tuning, which typically recovers 40–50% of the gap.

2.7.1 ACCURACY VS. EFFICIENCY TRADE-OFF ANALYSIS

The characterization of a 10-point drop as ”catastrophic” requires contextualization within the compression-accuracy Pareto frontier:

Interpretation:

- FuseGPT achieves an efficiency ratio of 0.132, comparable to state-of-the-art unstructured pruning (SparseGPT: 0.129) while offering hardware-friendly structured sparsity.
- The 1.33 \times speedup with 10.1-point drop is **substantially more favorable** than alternatives: LaCo’s 1.22 \times speedup costs 15.8 points, yielding a 50% worse efficiency ratio.
- Industry deployment studies (e.g., Google’s BERT pruning) suggest that efficiency ratios \geq 0.10 are generally acceptable for production use, placing FuseGPT well within practical bounds.

594
595
596
597
598
599
600
601
602
603
604
605
606
607
608
609
610
611
612
613
614
615
616
617
618
619
620
621
622
623
624
625
626
627
628
629
630
631
632
633
634
635
636
637
638
639
640
641
642
643
644
645
646
647

Table 20: Pareto analysis: accuracy degradation vs. speedup across methods (LLaMA-2-7B).

Method	Avg ZS Drop (points)↓	Speedup↑	Efficiency Ratio [†]	Acceptable ? (>0.10)
SliceGPT (25%)	-16.2	1.13×	0.070	× No
LaCo (25%)	-15.8	1.22×	0.077	× No
SLEB (25%)	-12.8	1.13×	0.088	× No
FuseGPT (25%)	-10.1	1.33×	0.132	✓ Yes
<i>Reference: unstructured pruning</i>				
SparseGPT (50% 2:4)	-8.5	1.10×	0.129	✓ Yes
Wanda (50% 2:4)	-9.2	1.10×	0.120	✓ Yes

[†]Efficiency Ratio = Speedup / ZS Drop; higher is better. Threshold of 0.10 based on industry practice.