

COBRA: Contribution-Based Bayesian Rank Allocation for Parameter-Efficient Fine-Tuning

Hongcheng Ding*, Xuanze Zhao*†, Liu Xuanhuang, Jing Jin, Shamsul Nahar Abdullah, Deshinta Arrova Dewi

Dongfang College, Zhejiang University of Finance & Economics · INTI International University · SEGi University

*Equal Contribution †Corresponding Author 20070032@zufedf.c.edu.cn, I24025877@sudent.newinti.edu.my



Abstract

Full fine-tuning of large language models (LLMs) incurs prohibitive computational and storage costs. Parameter-efficient fine-tuning (PEFT) addresses this limitation, with Low-Rank Adaptation (LoRA) gaining widespread adoption due to its simplicity and zero inference overhead. However, LoRA and its variants typically rely on uniform rank allocation or a single importance metric such as gradient magnitude or output sensitivity to guide rank distribution. This approach fails to recognize that gradient magnitude and output contribution are decoupled properties, leading to suboptimal allocation where critical layers are under-provisioned while less important ones waste capacity.

We introduce **COBRA** (Conductance-based Bayesian Rank Allocation), it operates in three core stages to a principled framework integrating dual importance factors for adaptive rank allocation in PEFT.

- S1 Layer Conductance:** Path-integral attribution for layer contribution $\phi^{(l)}$
- S2 Dual-Factor Aggregation:** Multiplicative fusion with gradient $\Gamma^{(l)} \rightarrow$ TA-LC prior $\pi^{(l)}$
- S3 Bayesian Allocation:** ELBO-driven variational rank optimization $\{r_l\}$

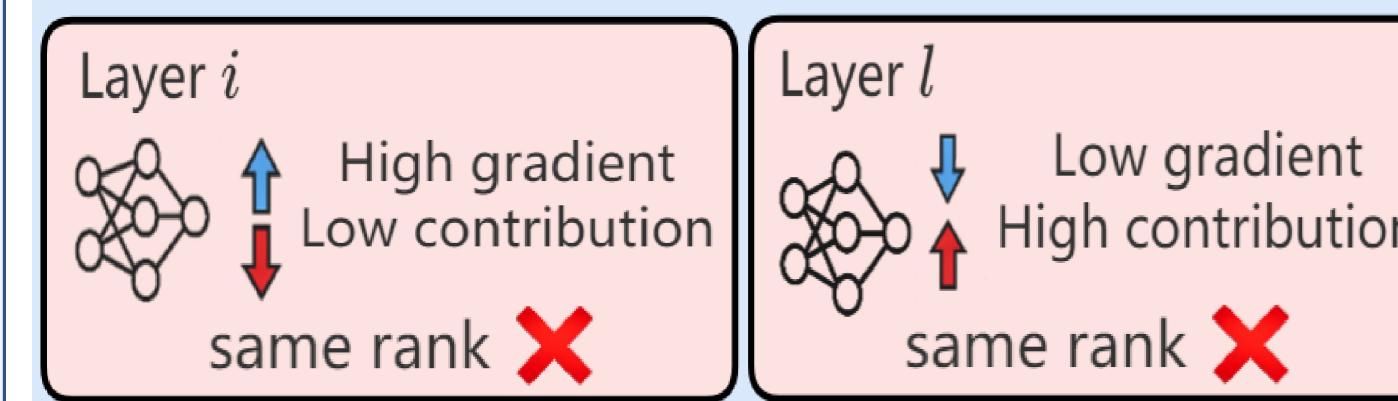
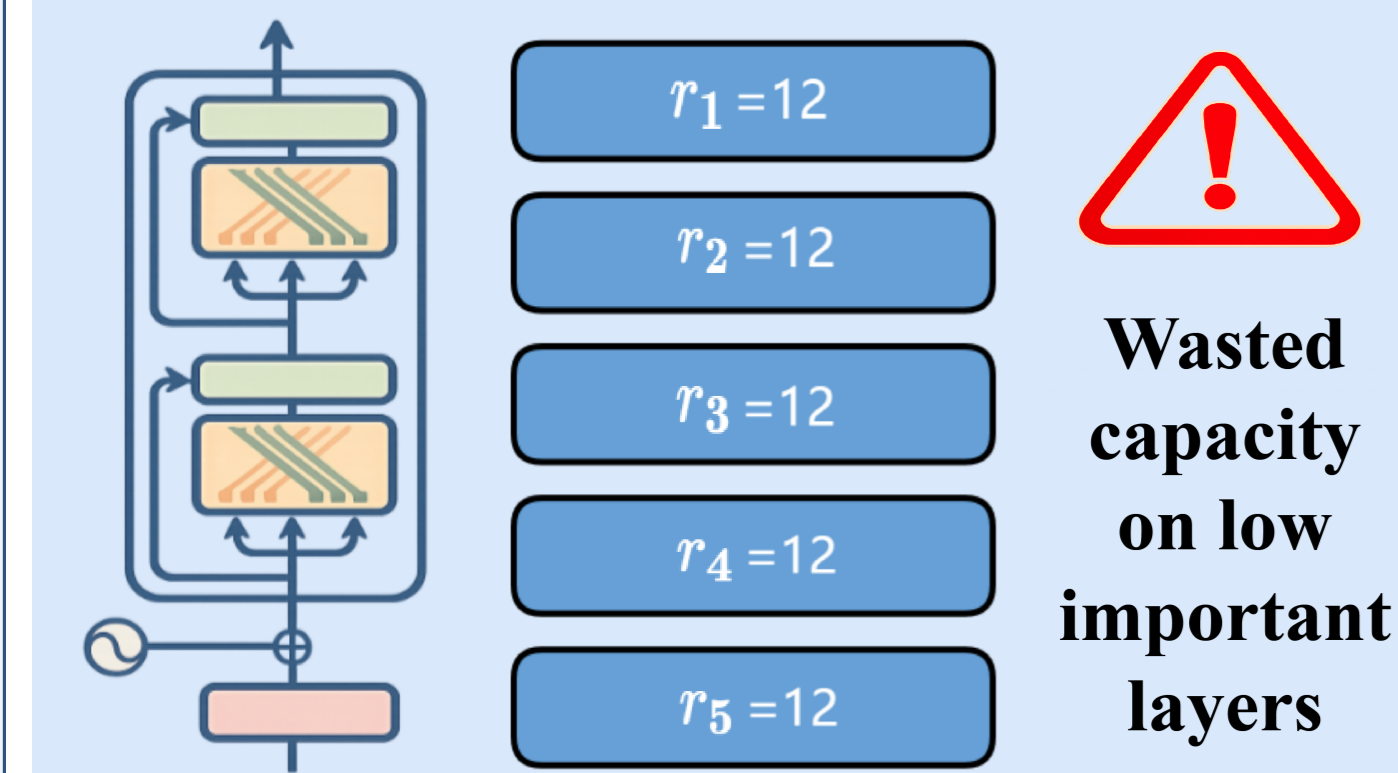
On GLUE, MMLU, GSM8K, HumanEval, and financial regression benchmarks, **COBRA** achieves **+1.6 pts** on GLUE, **+0.7%** over AdaLoRA, **+2.5%** over GoRA on LLaMA-7B, and **+6.6%** in high-rank regimes, under *identical* parameter budgets and with **zero inference overhead**.

- P1 Uniform Rank Allocation:** Static equal ranks neglect heterogeneous layer functional importance.
- P2 Single-Metric Deficiency:** Isolated metrics fail to capture decoupled layer contribution and adaptation demand.
- P3 Parameter Capacity Mismatch:** Heuristic rank allocation causes critical layer under-provisioning and parameter waste.

Uniform Rank Allocation

$$\Delta W^{(l)} = A^{(l)} \times B^{(l)}$$

(uniform rank r for all layers)



COBRA: Dual-Factor Adaptive Allocation

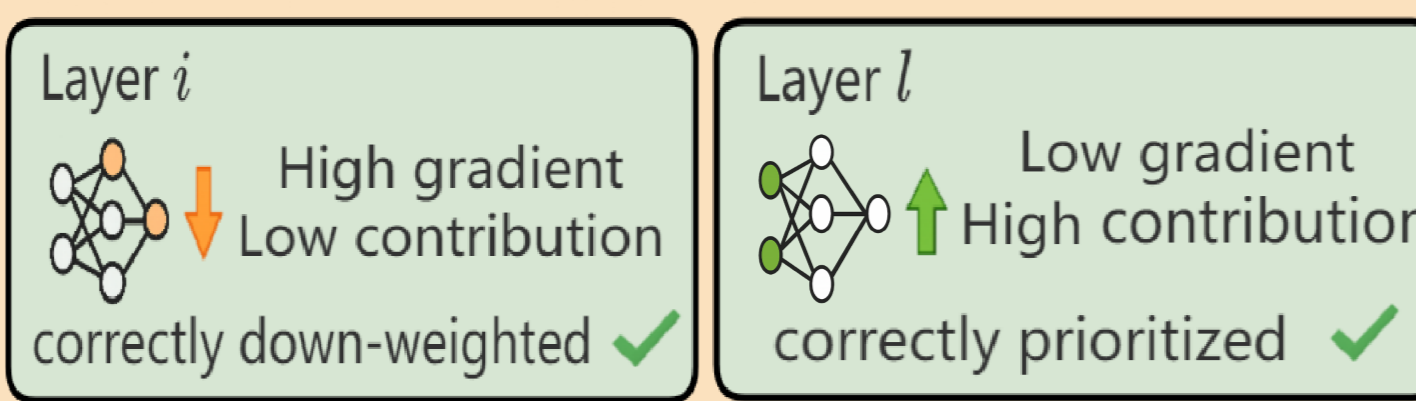
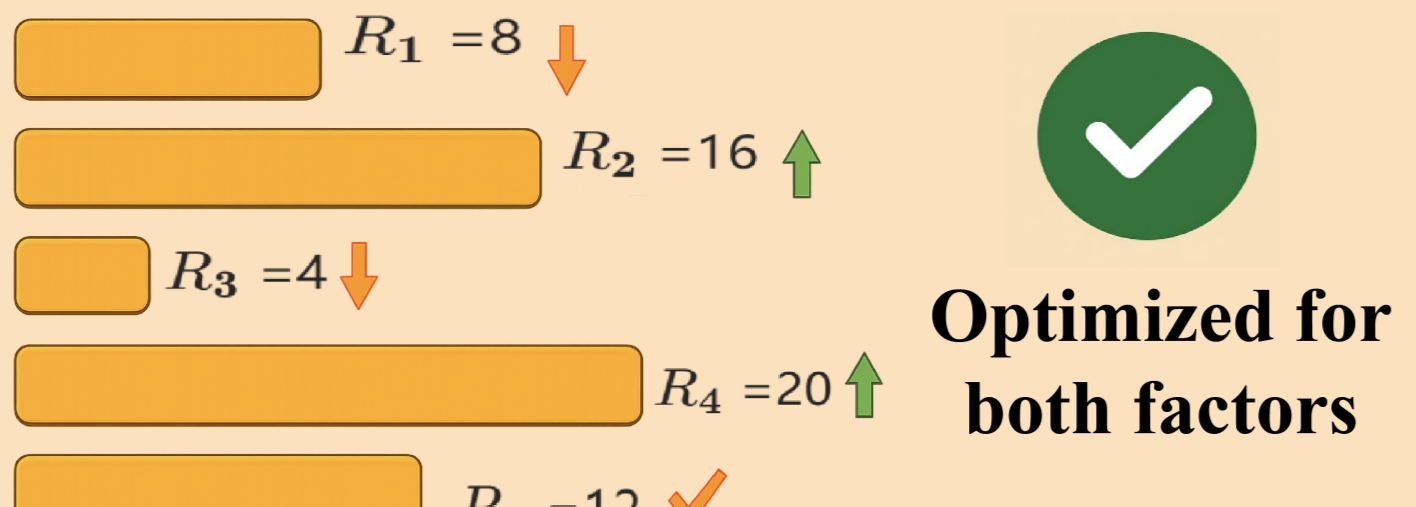
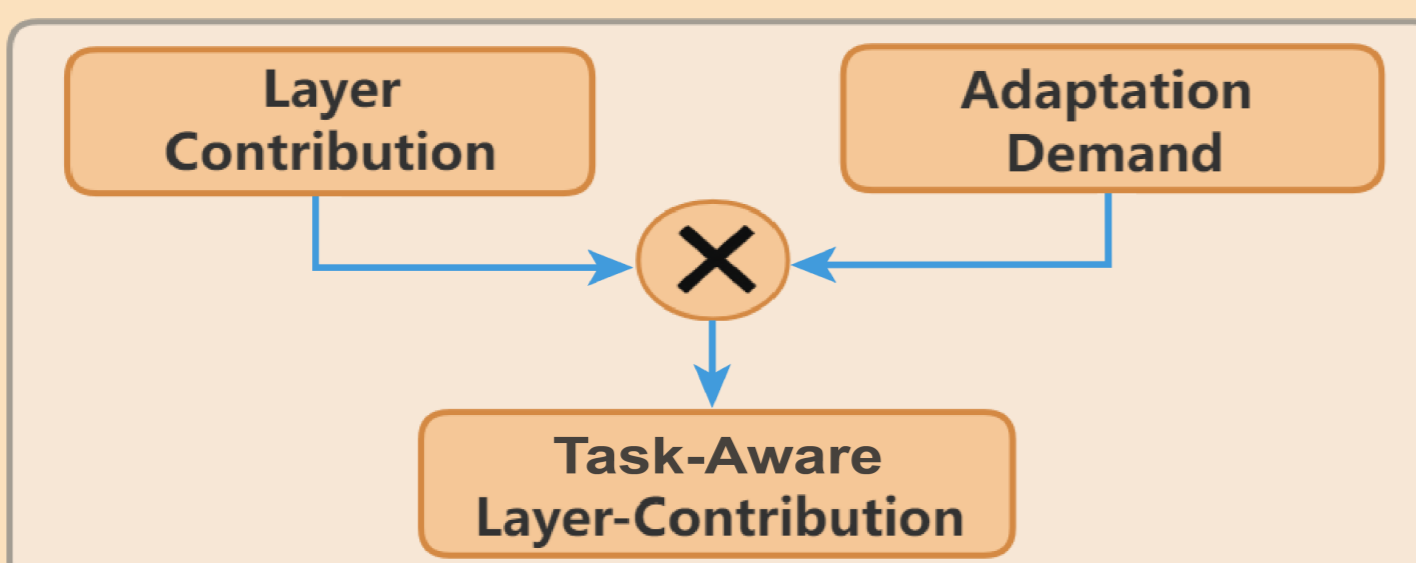


Figure 1: (Left) Uniform rank wastes capacity on low-contribution layers. (Right) COBRA's dual-factor adaptive allocation correctly prioritizes layers that are both *contributive* and *adaptation-demanding*.

Motivation: Two Critical Gaps

We probe **7 architectures** (encoder/decoder/enc-dec, 86M-7B params) and uncover two failure modes in single-factor rank allocation:

- Challenge #1: Layer Heterogeneity**
 - Gradient magnitude varies $9 \times -18 \times$ across layers
 - Output contribution varies $10 \times -19 \times$
 - Peaks diverge:** contribution peaks at *middle* layers (L13-16); gradient peaks at *top* layers (L21-24)
 - \Rightarrow Uniform allocation systematically wastes budget

- Challenge #2: Decoupled Characteristics**
 - Gradient \perp contribution: only *moderately* correlated ($\rho=0.66$ RoBERTa-large, $\rho=0.56$ LLaMA-2-7B)
 - Failure A:** high gradient + low contribution \rightarrow rank wasted on minimal-impact layers (GoRA's blind spot)
 - Failure B:** high contribution + moderate gradient \rightarrow critical layers under-provisioned (AdaLoRA's blind spot)
 - \Rightarrow *Neither metric alone suffices*

These gaps motivate a **dual-factor, Bayesian** framework that integrates contribution and adaptation demand into a single principled allocation rule.

Key Contributions

- C1. Decoupling Diagnosis:** First systematic evidence that gradient magnitude and output contribution are decoupled across 7 transformer architectures, exposing the root cause of single-factor allocation failure.
- C2. Interpretability-Guided PEFT:** A novel *Task-Adaptive Layer Conductance* (TA-LC) distribution combining path-integral attribution with adaptation demand, embedded as Bayesian prior in a variational rank-allocation objective.
- C3. Strong, Scalable Gains:** **+1.6** GLUE (RoBERTa-base), **+0.7%** over AdaLoRA, **+5.9%** over LoRA on LLaMA-7B, **+6.6%** in high-rank regimes; advantage *grows* with parameter budget.

Method: COBRA Framework

COBRA jointly leverages **two complementary signals** unified via Bayesian inference:

$$\phi^{(l)}: \text{Layer Contribution} \quad \Gamma^{(l)}: \text{Adaptation Demand} \quad \pi^{(l)}: \text{TA-LC Prior}$$

Stage 1: Layer Conductance Attribution Quantify exhaustive, non-redundant layer contribution via path-integral attribution:

$$A^{(l)}(x) = \sum_{j=1}^d (h_j^{(l)}(x) - h_j^{(l)}(x_0)) \cdot \int_0^1 \frac{\partial F(x_\alpha)}{\partial h_j^{(l)}} d\alpha, \quad \phi^{(l)} = \frac{1}{|\mathcal{D}_{\text{cal}}|} \sum_{x \in \mathcal{D}_{\text{cal}}} |A^{(l)}(x)| \quad (1)$$

- Calibration set:** $x_\alpha = x_0 + \alpha(x - x_0)$, 50 steps, $|\mathcal{D}_{\text{cal}}|=1000$
- Completeness:** $\sum_l A^{(l)}(x) = F(x) - F(x_0) \Rightarrow$ exhaustive & non-redundant
- One-time computation, amortized across all training runs

Stage 2: Dual-Factor Aggregation \rightarrow TA-LC Fuse layer contribution and adaptation demand multiplicatively:

$$\Gamma^{(l)} = \sqrt{\mathbb{E}_{(x,y) \sim \mathcal{D}_{\text{cal}}} [\|\nabla_{W_0^{(l)}} \mathcal{L}(x,y)\|_2^2]}, \quad \tilde{\phi}^{(l)} = \frac{\phi^{(l)}}{\sum_{l=1}^L \phi^{(l)}}, \quad \tilde{\Gamma}^{(l)} = \frac{\Gamma^{(l)}}{\sum_{l=1}^L \Gamma^{(l)}}, \quad \Psi^{(l)} = \tilde{\phi}^{(l)} \cdot \tilde{\Gamma}^{(l)}, \quad \pi^{(l)} = \frac{\Psi^{(l)}}{\sum_{l=1}^L \Psi^{(l)}} \quad (2)$$

- Multiplicative** (not additive): non-compensatory \rightarrow both factors must matter
- Down-weights layers that score high on *only one* factor (Failure A & B)
- Yields a normalized probability distribution over layers with $\sum_{l=1}^L \pi^{(l)} = 1$.

Stage 3: Bayesian Rank Allocation

$$\Delta W^{(l)} = \sum_{k=1}^{r_{\text{max}}} z_k^{(l)} \mathbf{a}_k^{(l)} (\mathbf{b}_k^{(l)})^\top, \quad \rho(z_k^{(l)}=1) = \rho^{(l)} = \sigma(\gamma(\pi^{(l)} - \bar{\pi})) \quad (3)$$

- Rank-level gates $z_k^{(l)} \in \{0, 1\}$; TA-LC-informed Bernoulli prior
- Continuous relaxation + Gumbel reparameterization for end-to-end gradients
- Top- R_{total} thresholding extracts deterministic $\{r_l\}$ under budget

Variational Objective (ELBO):

$$\mathcal{J}_{\text{ELBO}} = \mathbb{E}_q[\log p(\mathcal{D} | \theta, \tilde{\mathbf{z}})] - \beta \cdot D_{\text{KL}}(q(\tilde{\mathbf{z}}) \| p(\tilde{\mathbf{z}})) \quad (4)$$

$\theta = \{A^{(l)}, B^{(l)}\}$: LoRA params; $\phi = \{\alpha_k^{(l)}\}$: variational params; $\beta = 0.01, \gamma = 5.0, \lambda: 1.0 \rightarrow 0.1$ (annealed).

Key Design Advantages:

- Dual-factor:** fixes gradient-contribution decoupling
- Bayesian:** smooth optimization, no hard pruning instability
- Zero inference overhead:** same as standard LoRA
- Amortized calibration:** one-time cost reused across runs

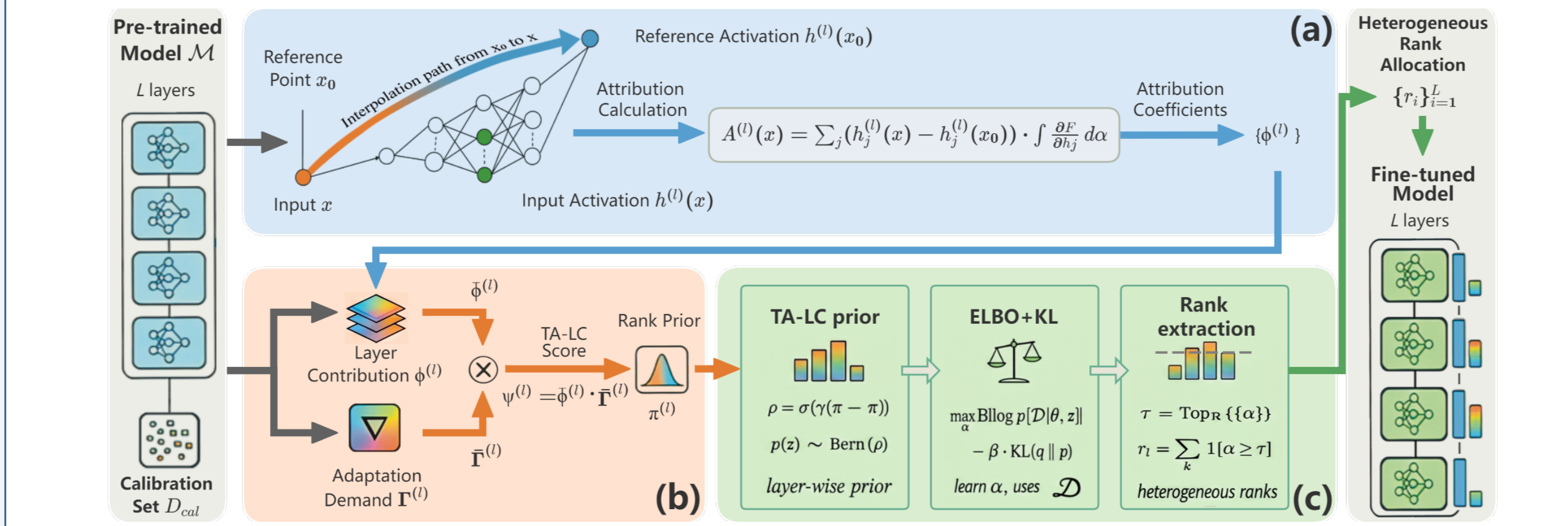


Figure 2: (a) Conductance attribution via interpolation path. (b) Multiplicative fusion of contribution & adaptation demand into TA-LC prior. (c) ELBO-based Bayesian gating with threshold extraction \rightarrow heterogeneous ranks $\{r_l\}_{l=1}^L$.

Training Protocol & Implementation

Algorithm 1 COBRA Procedure

Require: Pre-trained model \mathcal{M} , calibration set \mathcal{D}_{cal} , training set \mathcal{D} , budget R_{total} , max rank r_{max}

- Stage 1: Dual-Factor Estimation**
- Compute $\{\phi^{(l)}\}$ and $\{\Gamma^{(l)}\}$ via Eq. (1)-(2)
- Stage 2: Prior Construction**
- Compute $\{\pi^{(l)}\}$ and $\{\rho^{(l)}\}$ via Eq. (2)-(3)
- Stage 3: Bayesian Training**
- Initialize $A^{(l)} \sim \mathcal{N}$, $B^{(l)} \leftarrow 0$, $\alpha_k^{(l)} \leftarrow \rho^{(l)}$
- for each training iteration **do**
- Sample $\tilde{z}^{(l)}$ via Eq. (10)
- Compute $\mathcal{J}_{\text{ELBO}}$ via Eq. (4)
- Update θ, ϕ via gradient descent
- end for**
- Stage 4: Rank Extraction**
- $\tau \leftarrow \text{Top}_{R_{\text{total}}}(\{\alpha_k^{(l)}\}_{l,k})$
- for $l = 1$ to L **do**
- $r_l \leftarrow \sum_{k=1}^{r_{\text{max}}} \mathbf{1}[\alpha_k^{(l)} \geq \tau]$
- end for**

Ensure: Fine-tuned model with $\{r_l\}_{l=1}^L$

Hyperparameters & Setup:

- Adapters on *all* linear layers (attention + FFN)
- AdamW, lr $\in \{1, 2, 5\} \times 10^{-4}$, batch 16/32, 6% warmup, linear decay
- Rank budgets: $R_{\text{total}} \in \{192, 384, 768\}$ (NLU) & 9K, 12K (finance)
- Hardware: NVIDIA A100 80GB, FP16 activation / FP32 gradient

Evidence: Gradient \neq Contribution

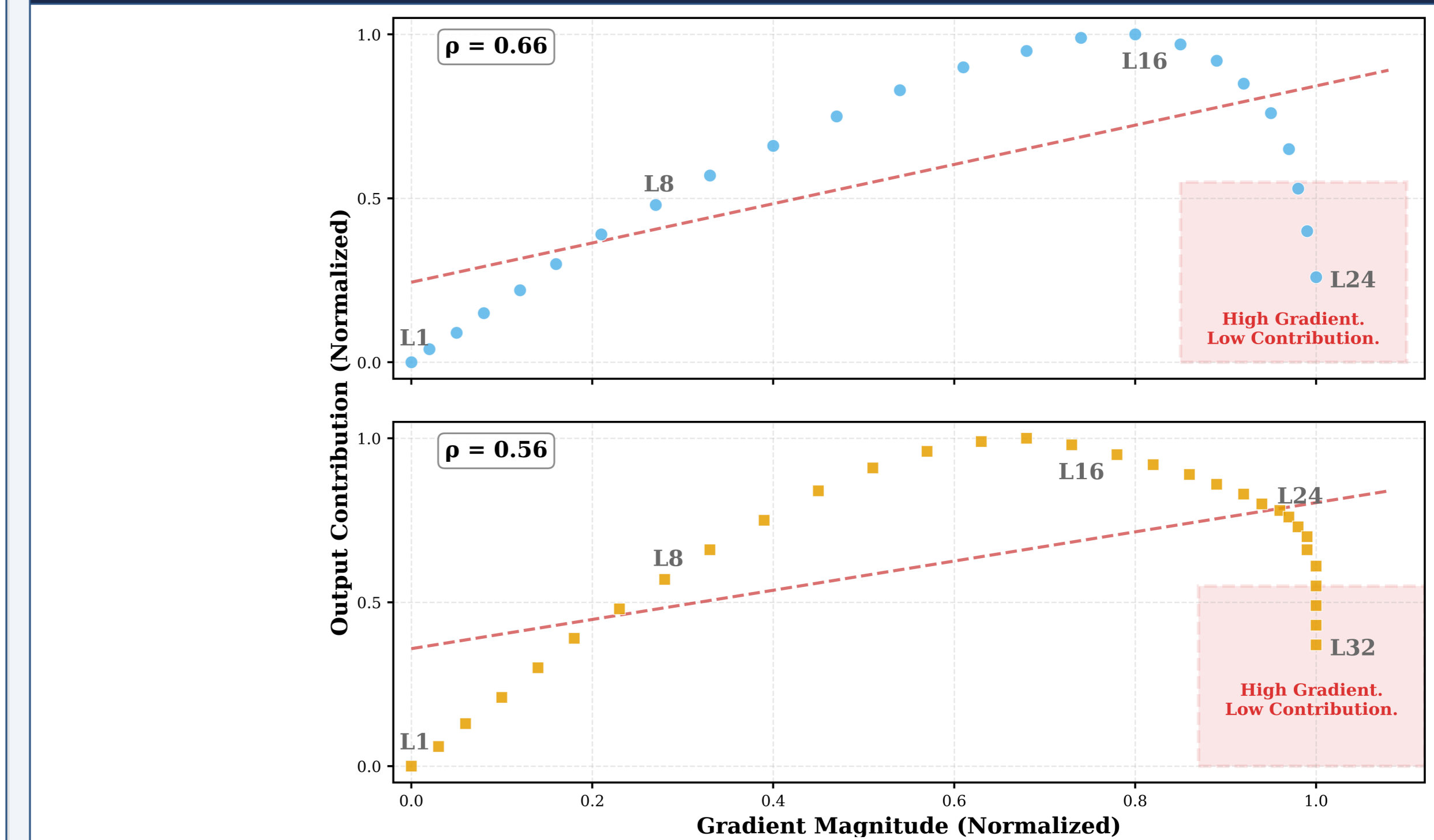


Figure 3: Gradient vs. contribution for RoBERTa-large ($\rho=0.66$) & LLaMA-2-7B ($\rho=0.56$). Red zones: high-gradient, low-contribution layers that gradient-driven methods over-provision.

- Failure A:** L24/L32 sit in the *red zone* \rightarrow max gradient, minor contribution (GoRA over-allocates)
- Failure B:** Mid-layers (L8-L16) above trend \rightarrow high contribution, moderate gradient (under-allocated)
- COBRA's TA-LC** ($\Phi \times \Gamma$) correctly suppresses red zone & amplifies mid-layers

Experiments & Main Results

NLU & generalization \rightarrow GLUE (RoBERTa-b / DeBERTaV3-b) & CSQA (Gemma-2-2B), identical budgets; **COBRA wins all three:**

Table: Low-rank benchmarks vs. adaptive baselines (\uparrow better). Uniform LoRA shaded gray; COBRA in pink.

Method	Type	RoBERTa-b	DeBERTaV3-b	Gemma-2-2B
		GLUE $r=8$	GLUE $r=4$	CSQA $r=8$
LoRA	Uniform	87.0	88.7	74.2
DyLoRA	Dyn. range	87.1	88.4	74.5
AutoLoRA	NAS/Meta	87.5	89.0	75.6
DoRA-dyn	Dyn. dist.	87.9	89.3	75.9
GoRA	Gradient	87.6	88.8	75.1
COBRA	Dual-factor	88.6	89.9	76.1

LLM scaling (vs. gradient-driven GoRA, LLaMA-7B, $r_0=8$):

- MMLU **48.6** (+1.5), GSM8K **45.9** (+2.1), HumanEval **52.7** (+3.7)
- Relative gain over LoRA **+5.9%** vs. GoRA's +3.4%

High-rank regression financial sentiment, 3 datasets:

Table: MSE \downarrow across rank budgets. LoRA gray (baseline); COBRA pink. Gains grow with budget (avg. **+6.6%**).

Method	Config.	NEU		FXE		INV	
		MSE	MAE	MSE	MAE	MSE	MAE
Full FT	-	0.0189	0.0981	0.0202	0.1014	0.0183	0.0959
LoRA	$r=384$	0.0189	0.0977	0.0205	0.1021	0.0184	0.0961
	$r=512$	0.0186	0.0958	0.0199	0.1007	0.0183	0.0955
AdaLoRA	640 \rightarrow 512	0.0218	0.1052	0.0231	0.1079	0.0196	0.1005
COBRA	$r_{\text{avg}}=384$	0.0177	0.0940	0.0192	0.0978	0.0172	0.0917
	$r_{\text{avg}}=512$	0.0174	0.0931	0.0185	0.0958	0.0170	0.0908

Ablation Study (DeBERTaV3-base)

Factor Aggregation Variant	Acc.	Prior γ (inverted-U)		KL Weight β	
		Variant	Acc.	Variant	Acc.
Φ/Γ only	88.8	$\gamma=0$ (none)	88.5	$\beta=0$	88.8
$\Phi+\Gamma$ (add.)	89.2	$\gamma=50$ (over)	88.3	$\beta=1.0$ (over)	88.2
$\Phi \times \Gamma$ (ours)	89.9	$\gamma=5.0$ (ours)	89.9	$\beta=0.01$ (ours)	89.9

- Multiplicative** fusion is non-compensatory: a layer must score high on *both* factors (+0.7 over additive).
- Interpretability-guided prior is essential: weak \rightarrow unguided, strong \rightarrow over-constrained.

Conclusion & Takeaways

COBRA resolves the gradient-contribution decoupling that limits prior PEFT by unifying **interpretability** (path-integral conductance) with **adaptation demand** under a **Bayesian** objective.

- Dual-factor** allocation fixes single-metric blind spots (GoRA & AdaLoRA).
- Bayesian relaxation** gives stable optimization, no hard-pruning instability.
- Gains **scale with budget** (+1.6 GLUE \rightarrow +6.6% high-rank); **zero inference overhead**.
- Model-agnostic:** encoders, decoders & encoder-decoders.

Future: Vision-language PEFT · Causal attribution · Auto-budget allocation.