

POST-GATING BIAS: RESTORING AFFINE FREEDOM IN TRANSFORMER MLPs

Anonymous authors

Paper under double-blind review

ABSTRACT

Modern transformers often omit additive biases because normalization and attention preserve constant offsets that downstream linear maps can remap. The gated MLP (SwiGLU) is a notable exception: the elementwise product with a nonlinear gate destroys such offsets, removing affine freedom after the nonlinearity. We examine a simple modification—Post-Gating Bias (PGB), an additive term applied after the gated product and before the down-projection. PGB restores this degree of freedom with negligible computational cost. Our working hypothesis is that PGB mitigates training noise from dropout or from additive stochastic regularization (e.g., in VAEs) by shifting activation boundaries in a controlled way, thereby softening sharp transitions that otherwise amplify perturbations. We observe stability gains at higher learning rates and some robustness improvements in a ViT-VAE setting. We also show other settings where the effect is minimal. We present these observations to clarify where biases are redundant in transformer blocks and where multiplicative gating makes them potentially useful. Finally, we report a controlled study varying dropout and latent noise across multiple seeds to test this hypothesis.

1 INTRODUCTION

Transformer architectures rely heavily on normalization layers to stabilize training. Because layer normalization (LN) and RMS normalization (RMSNorm) both incorporate additive biases, it is often assumed that explicit bias terms in intermediate linear transformations are redundant: any constant offset can be reproduced elsewhere in the block at negligible cost to expressivity. This assumption has motivated widespread omission of biases in recent transformer designs.

We revisit this assumption in the context of gated multilayer perceptrons (MLPs), the feed-forward design popularized by Shazeer (2020) and now used in nearly all state-of-the-art large language models. SwiGLU, in particular, is the standard choice in PaLM (Chowdhery et al., 2022), LLaMA (Touvron et al., 2023), Phi-3 (Abdin et al., 2024), and Gemma-3 (Team, 2024). While the redundancy of biases generally holds throughout most of the transformer block, it breaks down at a critical point: the final down-projection after elementwise gating. At this stage, signal modulation by the gate prevents constant offsets from being trivially reconstructed, so omitting a bias can interact with dropout or latent noise in ways that destabilize training.

To address this, we introduce a *post-gating bias* (PGB), added after the gated multiplication but before the down-projection:

$$U = XW_u, \tag{1}$$

$$G = XW_g, \tag{2}$$

$$Y = \text{dropout}(U) \odot \text{silu}(G) + \mathbf{b}, \tag{3}$$

$$Z = YW_d. \tag{4}$$

This modification incurs virtually no additional parameter or memory cost, as the bias can be applied in place, yet it can smooth optimization when stochastic noise is present.

Contributions. Our work makes three contributions: (i) we revisit the assumption that biases in transformer blocks are redundant, providing an analysis of when biases can be absorbed into weight

054 matrices and showing that this equivalence fails after multiplicative gating; (ii) we propose the *post-*
 055 *gating bias* (PGB), a negligible-cost addition that restores affine freedom at the down-projection
 056 stage; (iii) we present preliminary experiments in three settings—a ViT-VAE on CelebA (Lee et al.,
 057 2020; Chang et al., 2022), modified Phi-3 language models (Abdin et al., 2024), and gated MLP
 058 classifiers—which suggest that PGB can improve stability in noisy latent settings and may enhance
 059 adversarial robustness, while showing little to no effect in large-scale autoregressive pretraining.

060 These findings indicate that even apparently redundant architectural elements can play subtle roles
 061 under stochastic training conditions, motivating further exploration in settings where noise and
 062 dropout strongly shape optimization dynamics.

063 2 BACKGROUND AND RELATED WORK

064
 065
 066
 067
 068 Early work on gated architectures demonstrated the effectiveness of gating in recurrent models
 069 (Jozefowicz et al., 2015; 2016), inspiring Dauphin et al. (2017) to introduce gated linear units
 070 (GLUs) in convolutional networks. Their formulation $(\mathbf{X}\mathbf{W} + \mathbf{b}) \odot \sigma(\mathbf{X}\mathbf{V} + \mathbf{c})$ included bi-
 071 ases on all affine maps, and gating occurred in the same dimensionality as the input, with residual
 072 connections across convolutional blocks.

073 Building on this idea, Shazeer (2020) introduced GLU variants for transformer feed-forward net-
 074 works, including ReGLU, GEGLU, and SwiGLU. Their design uses an explicit up-gate-down pro-
 075 jection structure, but notably omits all bias terms in the MLP layers, including the down-projection.
 076 This bias-free formulation has since influenced large-scale LLMs, where biases in dense layers are
 077 often removed entirely. Subsequent works such as So et al. (2021) integrated gating variants into
 078 architecture search, and Ramesh & Ramkumar (2023) extended gated activations to attention blocks
 079 in vision transformers. While these contributions reinforce the importance of gating in transformer
 080 design, none have revisited the role of bias placement relative to dropout and gating.

081 In parallel, several lines of work have explored transformer-based VAEs. Conditional generation
 082 models such as T-CVAE (Wang et al., 2019) and TRACE (Hu et al., 2022) adapt variational bottle-
 083 necks for diverse text generation. Other approaches embed latent variables directly in transformer
 084 attention layers via nonparametric variational information bottlenecks (Henderson & Fehr, 2022),
 085 or construct structured latent spaces such as graph-induced syntactic-semantic embeddings (Zhang
 086 et al., 2023). Beyond NLP, transformer-VAE hybrids have been proposed for manifold-aware vision
 087 representation learning (Shamsolmoali et al., 2023) and reduced-order modeling in fluid dynamics
 088 (Solera-Rico et al., 2023). Collectively, these works show the versatility of combining transformers
 089 with variational inference, but they do not address how architectural details like bias and dropout
 090 placement in gated MLP blocks affect optimization stability.

091 Our experimental study confirms the prevailing view in language modeling—biases after gating do
 092 not improve perplexity in standard transformer LMs—but we also find evidence in autoencoding
 093 settings that post-gating biases can stabilize training under latent noise. This suggests their utility
 094 may depend not on expressivity but on the interaction between dropout, gating, and the statistical
 095 structure of the task.

096 3 THEORETICAL ANALYSIS

097
 098
 099
 100
 101 The work above suggests that bias placement does not affect model expressivity in most parts of
 102 the transformer block, but that gating combined with dropout or latent noise creates an exception:
 103 offsets routed through the multiplicative interaction become unstable.

104 To make this precise, we analyze when biases can be reconstructed from weights alone and when
 105 such reconstruction becomes ill-conditioned. We show that (i) constant offsets can be re-encoded
 106 throughout attention and other linear maps, (ii) this equivalence breaks down after multiplicative gat-
 107 ing, where the mean direction itself carries variance, and (iii) a post-gating bias restores conditioning
 by centering the regressors seen by the down-projection. Proofs are deferred to Appendix A.

3.1 OUTPUT-BIAS RECONSTRUCTABILITY

We compare two realizations of the same affine map,

$$\mathbf{y}_1 = \mathbf{W}(\mathbf{x} + \mathbf{b}_{\text{in}}) + \mathbf{b}_{\text{out}} \quad \text{vs.} \quad \mathbf{y}_2 = \hat{\mathbf{W}}(\mathbf{x} + \mathbf{b}_{\text{in}}), \quad \hat{\mathbf{W}} = \mathbf{W} + \Delta,$$

and ask for the *smallest* perturbation Δ that reconciles the missing intercept.

Exact-fit with minimal change. We enforce agreement on the baseline direction and minimize the change:

$$\min_{\Delta} \|\Delta\|_F \quad \text{s.t.} \quad \Delta \mathbf{b}_{\text{in}} = \mathbf{b}_{\text{out}}. \quad (5)$$

Lemma 1 (Minimal-change exact-fit is rank-1). *The unique minimizer of equation 5 is*

$$\Delta^* = \frac{\mathbf{b}_{\text{out}} \mathbf{b}_{\text{in}}^\top}{\|\mathbf{b}_{\text{in}}\|_2^2}, \quad \hat{\mathbf{W}} = \mathbf{W} + \Delta^*.$$

This perturbation exactly recovers the output bias and leaves all components of \mathbf{x} orthogonal to \mathbf{b}_{in} unaffected. Residual discrepancies are confined to the \mathbf{b}_{in} direction and $\|\Delta\|_F$ is minimal.

Operator-norm remark. Among all Δ satisfying $\Delta \mathbf{b}_{\text{in}} = \mathbf{b}_{\text{out}}$, the rank-1 map in Lemma 1 minimizes all spectral norms. Any action off $\text{span}(\mathbf{b}_{\text{in}})$ would increase the operator norm and the worst-case discrepancy.

The construction above shows that weight-only fixes *exist*. However, first-order training does not target equation 5 explicitly. In practice, SGD may not discover Δ^* , and when dropout or latent noise is present, routing the intercept through $(\mathbf{x} + \mathbf{b}_{\text{in}})$ couples it to high-variance, intermittently masked directions, degrading conditioning and stability. A post-gating bias supplies the intercept directly, avoiding this mean-through-weights pathway.

Remark (Minimum-MSE perturbation). Let \mathbf{x} have mean 0 and covariance $\Sigma \succeq 0$. The solution of

$$\min_{\Delta} \mathbb{E} \left\| (\mathbf{W}(\mathbf{x} + \mathbf{b}_{\text{in}}) + \mathbf{b}_{\text{out}}) - (\mathbf{W} + \Delta)(\mathbf{x} + \mathbf{b}_{\text{in}}) \right\|_2^2$$

is

$$\Delta^* = \mathbf{b}_{\text{out}} \mathbf{b}_{\text{in}}^\top (\Sigma + \mathbf{b}_{\text{in}} \mathbf{b}_{\text{in}}^\top)^{-1}, \quad \mathbb{E} \|\mathbf{y}_1 - \mathbf{y}_2\|_2^2 \Big|_{\min} = \|\mathbf{b}_{\text{out}}\|_2^2 \left(1 - \mathbf{b}_{\text{in}}^\top (\Sigma + \mathbf{b}_{\text{in}} \mathbf{b}_{\text{in}}^\top)^{-1} \mathbf{b}_{\text{in}} \right).$$

As $\Sigma \rightarrow 0$, this reduces continuously to the rank-1 exact-fit in Lemma 1. Despite existence, SGD has no reason to recover this projection; hence we defer the proof to Appendix A and emphasize the stability benefits of an explicit post-gating bias in the main text.

Reconstructability through attention and the output projection. Let $\hat{\mathbf{X}} = \text{RMSNorm}(\mathbf{X}) + \mathbf{b}$ and $[\mathbf{Q}, \mathbf{K}, \mathbf{V}] = \hat{\mathbf{X}}[\mathbf{W}_q, \mathbf{W}_k, \mathbf{W}_v]$. Any constant offset \mathbf{b} can be linearly propagated into $\mathbf{Q}, \mathbf{K}, \mathbf{V}$. Crucially, if \mathbf{V} carries a constant offset per token, i.e., $\mathbf{V} = \hat{\mathbf{V}} + \mathbf{1} \mathbf{c}^\top$ for some $\mathbf{c} \in \mathbb{R}^{d_v}$ (with $\mathbf{1}$ the all-ones column), then with standard attention

$$\mathbf{P} = \text{softmax}\left(\frac{\mathbf{Q}\mathbf{K}^\top}{\sqrt{d_k}}\right) \quad (\text{row-stochastic: } \mathbf{P}\mathbf{1} = \mathbf{1}),$$

the attention output is

$$\mathbf{A} = \mathbf{P}\mathbf{V} = \mathbf{P}\hat{\mathbf{V}} + \mathbf{P}(\mathbf{1} \mathbf{c}^\top) = \mathbf{P}\hat{\mathbf{V}} + \mathbf{1} \mathbf{c}^\top.$$

Thus the attention block *exactly* transmits one full bias vector per row via \mathbf{V} because the softmax weights are row-normalized. Consequently, the final output projection

$$\mathbf{Z} = \mathbf{A}\mathbf{W}_o = (\mathbf{P}\hat{\mathbf{V}})\mathbf{W}_o + \mathbf{1}(\mathbf{c}^\top \mathbf{W}_o)$$

can also reconstruct a constant offset if needed. This reinforces that biases are readily recoverable *throughout* attention; the non-reconstructability we highlight pertains specifically to the multiplicative gating stage prior to the down-projection in the MLP.

3.2 THE GATING EXCEPTION

The expressivity equivalence breaks down after multiplicative gating. In the MLP we write

$$U = \mathbf{X}W_u, \quad G = \mathbf{X}W_g, \quad S = U \odot \text{SiLU}(G), \quad Y = S \odot M + \mathbf{b}, \quad Z = YW_d,$$

with an optional dropout mask M and optional post-gating bias \mathbf{b} . If $\mathbf{b} = 0$, the down-projection W_d is the only route to produce offsets. Unlike attention, the multiplicative interaction $U \odot \text{SiLU}(G)$ destroys constant offsets; routing an offset through S couples it to high-variance, intermittently masked directions.

Ill-conditioning of weight-only reconstruction. The reasoning of Lemma 1 and the minimum-MSE remark applies not only to strict constant offsets but also to reconstruction of the *expected* output. In this case, the minimal-change solution is again rank-1, aligned with the mean of the post-gating features. However, variation along the mean direction is typically larger than in other coordinates, since when $\text{SiLU}(G)$ vanishes for a fraction of cases, the surviving coordinates both carry the expectation and accumulate higher variance. In high-dimensional latent spaces, information can ordinarily be encoded orthogonally to a bias direction, but under multiplicative gating the nonzero expectation itself *is* the signal, leaving no orthogonal subspace to carry offsets. Thus the rank-1 correction introduces a top-of-spectrum mode precisely in a high-variance direction, inflating the condition number of W_d and making outputs acutely sensitive to small perturbations. This explains why weight-only recovery of intercepts is particularly ill-conditioned in gated MLPs.

Proposition 1 (Intercept removal improves conditioning). *Let $\varphi \in \mathbb{R}^p$ denote the vectorized post-gating features (e.g., a row of $S \odot M$), with mean $\boldsymbol{\mu} = \mathbb{E}[\varphi]$ and covariance $\boldsymbol{\Sigma} = \text{Cov}(\varphi)$. Under squared loss, the (block) Hessian with respect to W_d without an explicit bias term is*

$$\mathbf{H}_{\text{no-bias}} = \mathbb{E}[\varphi\varphi^\top] = \boldsymbol{\Sigma} + \boldsymbol{\mu}\boldsymbol{\mu}^\top,$$

while with a trainable post-gating bias \mathbf{b} that absorbs the mean offset, the effective Hessian governing W_d is

$$\mathbf{H}_{\text{pgb}} = \text{Cov}(\varphi) = \boldsymbol{\Sigma}.$$

Since $\mathbf{H}_{\text{no-bias}} = \boldsymbol{\Sigma} + \boldsymbol{\mu}\boldsymbol{\mu}^\top$ adds a rank-1 positive semidefinite term, its largest eigenvalue satisfies

$$\lambda_{\max}(\mathbf{H}_{\text{no-bias}}) \geq \lambda_{\max}(\boldsymbol{\Sigma}) + (\mathbf{v}^\top \boldsymbol{\mu})^2$$

where \mathbf{v} is the normalized top-eigenvector of $\boldsymbol{\Sigma}$; consequently

$$\kappa(\mathbf{H}_{\text{no-bias}}) \geq \kappa(\boldsymbol{\Sigma}) \quad \text{with strict inequality when } \boldsymbol{\mu} \notin \ker(\boldsymbol{\Sigma}),$$

so centering via an explicit intercept strictly improves (or leaves unchanged) the condition number.

Implication. Without PGB, W_d must learn using the uncentered design matrix $\mathbb{E}[\varphi\varphi^\top] = \boldsymbol{\Sigma} + \boldsymbol{\mu}\boldsymbol{\mu}^\top$, whose additional rank-1 term both enlarges the top eigenvalue and—under dropout/latent noise—fluctuates across steps, creating a moving ill-conditioned direction. With PGB, the intercept is parameterized directly; gradients for W_d are driven by centered features with Hessian $\boldsymbol{\Sigma}$, yielding better conditioning and lower gradient variance.

Variance control via post-gating bias. Let T denote the (per-token) training target in the decoder’s space. Under a linear readout $Z = YW_d = (S \odot M)W_d + \mathbf{b}W_d$, the optimal intercept for MSE satisfies

$$\mathbf{b}^* W_d = \mathbb{E}[T] - \mathbb{E}[S \odot M] W_d.$$

At this solution, the residual driving updates of W_d are

$$\mathbf{R} = T - \mathbb{E}[T] - ((S \odot M) - \mathbb{E}[S \odot M]) W_d,$$

so the stochastic gradient w.r.t. W_d depends on the *centered* post-gating features $(S \odot M) - \mathbb{E}[S \odot M]$. By the matrix variance identity $\mathbb{E}[\varphi\varphi^\top] = \text{Cov}(\varphi) + \boldsymbol{\mu}\boldsymbol{\mu}^\top$, removing the mean eliminates the rank-1 $\boldsymbol{\mu}\boldsymbol{\mu}^\top$ contribution that otherwise dominates the spectrum and couples to noise-induced shifts of $\boldsymbol{\mu}$. Consequently, with PGB the effective gradient covariance is reduced from $\mathbb{E}[\varphi\varphi^\top]$ to $\text{Cov}(\varphi)$, improving optimization stability.

Remark (Centering upstream reduces second moments). Suppose \mathbf{X} has tokenwise expectations collected in rows of $\bar{\mathbf{X}}$ and let $\bar{\boldsymbol{\mu}}$ be the average row. Replacing \mathbf{X} by $\mathbf{X} - \mathbf{1}\bar{\boldsymbol{\mu}}^\top$ strictly reduces $\mathbb{E}\|\mathbf{X}\|_F^2$ by $\|\bar{\boldsymbol{\mu}}\|_2^2$ per token. If $\mathbf{U} = \mathbf{X}\mathbf{W}_u$ and $\mathbf{G} = \mathbf{X}\mathbf{W}_g$, this centering reduces $\mathbb{E}\|\mathbf{U}\|_F^2$ and $\mathbb{E}\|\mathbf{G}\|_F^2$ by $\|\bar{\boldsymbol{\mu}}\mathbf{W}_u\|_2^2$ and $\|\bar{\boldsymbol{\mu}}\mathbf{W}_g\|_2^2$, respectively. Since SiLU is 1-Lipschitz and monotone, a first-order bound gives

$$\mathbb{E}\|\mathbf{S}\|_F^2 = \mathbb{E}\|\mathbf{U} \odot \text{SiLU}(\mathbf{G})\|_F^2 \leq \mathbb{E}\|\mathbf{U}\|_F^2 \mathbb{E}\|\text{SiLU}(\mathbf{G})\|_\infty^2,$$

so upstream mean-removal lowers a surrogate bound on the second moment of \mathbf{S} . Trainable PGB lets the optimizer realize this reparametrization *in situ* by assigning baselines to \mathbf{b} and reserving \mathbf{U}, \mathbf{G} for deviations, thereby reducing the variance of \mathbf{Y} and improving conditioning of \mathbf{W}_d .

Stochastic perturbations from dropout. The analysis above shows that centering reduces the second moment of post-gating features in expectation. However, dropout (or other stochastic regularizers) reintroduces noisy offsets: even when $\mathbb{E}[\mathbf{S}] \approx 0$, the masked signal $\mathbf{M} \odot \mathbf{S}$ has expectation $p\boldsymbol{\mu}$ with variance inflated by $p(1-p)$. Thus, beyond deterministic centering, a trainable post-gating bias can serve as an *adaptive baseline* that cancels these noisy expectations. The following lemma formalizes this variance-reduction role.

3.3 VARIANCE REDUCTION UNDER DROPOUT

For sample i , the gradient contribution to \mathbf{W}_d is

$$\frac{\partial \ell_i}{\partial \mathbf{W}_d} = \mathbf{y}_i \boldsymbol{\gamma}_i^\top, \quad \mathbf{y}_i = \mathbf{M}_i \odot \mathbf{S}_i + \mathbf{b},$$

where \mathbf{M}_i is a dropout mask, $\mathbf{S}_i = \mathbf{U}_i \odot \text{SiLU}(\mathbf{G}_i)$, $\boldsymbol{\gamma}_i$ is the loss gradient, and \mathbf{b} is the post-gating bias. Row-wise,

$$\left(\frac{\partial \ell_i}{\partial \mathbf{W}_d} \right)_{j:} = (\mathbf{M}_{ij} \mathbf{S}_{ij} + \mathbf{b}_j) \boldsymbol{\gamma}_i^\top.$$

Dropout intermittently zeros \mathbf{S}_{ij} , while a nonzero \mathbf{b}_j supplies a stable baseline across all samples.

Lemma 2 (Baseline minimizes gradient variance). *Assume $\mathbb{E}[\mathbf{S}_{ij}] = \boldsymbol{\mu}_j$, $\text{Var}(\mathbf{S}_{ij}) = \sigma_j^2$, independence of \mathbf{M}_{ij} and $\boldsymbol{\gamma}_i$, and $\mathbb{E}[\boldsymbol{\gamma}_i] = 0$. Then for any test direction \mathbf{u} ,*

$$\text{Var}\left(\mathbf{u}^\top \left(\frac{\partial \ell_i}{\partial \mathbf{W}_d} \right)_{j:}^\top\right) = \mathbb{E}[(\mathbf{u}^\top \boldsymbol{\gamma}_i)^2] \left(p\sigma_j^2 + (p\boldsymbol{\mu}_j + \mathbf{b}_j)^2 \right),$$

which is minimized uniquely at $\mathbf{b}_j^* = -p\boldsymbol{\mu}_j$.

Interpretation. The post-gating bias can cancel the masked mean $p\boldsymbol{\mu}_j$, leaving only the irreducible term $p\sigma_j^2$. Thus PGB centers each regressor and strictly reduces gradient variance. Without PGB, the mean must be encoded through masked features, inflating variance and worsening conditioning.

4 EXPERIMENTS AND RESULTS

To evaluate the effectiveness of post-gating bias (PGB), we tested the central hypothesis across three representative settings: (i) a Vision Transformer VAE on CelebA, where stochastic latent encodings create noisy training signals, (ii) small-scale language model pretraining on OpenWebText, where autoregressive objectives dominate, and (iii) stacked gated MLP classifiers on MNIST Xiao et al. (2017) and CIFAR-10/CIFAR-100 Krizhevsky (2009), where robustness to adversarial perturbations can be directly assessed. Together, these experiments probe whether the benefits of PGB are specific to noisy latent autoencoders or extend to other architectures.

Experimental setting. For CelebA we use the standard training/validation/test splits. Our ViT-VAE employs 16×16 patches, an embedding dimension of 576, 12 attention heads, and 8 encoder and 8 decoder transformer blocks. Models are trained with Adam at a learning rate of 10^{-4} using a cosine schedule down to 10^{-5} , with dropout fixed at 0.1. The loss combines RMSE reconstruction error with a latent length penalty, and Gaussian noise is added to latent codes during training.

For language modeling we pretrain Phi3-mini and Gemma3 from scratch on OpenWebText, again using standard splits. Each model is trained with Adam at 10^{-4} with the same cosine schedule,

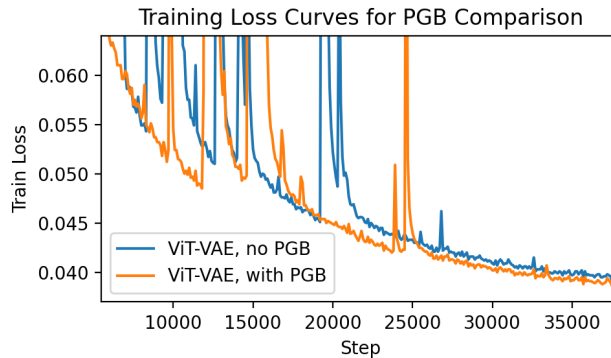


Figure 1: Training loss comparisons on a transformer-based VAE for CelebA. PGB shortens loss spikes and leads to longer stable regions of descent, with a slight asymptotic improvement.

and dropout probability 0.1. We evaluate test loss under four bias configurations: no bias, PGB only, RMSNorm bias only, and both biases. For Gemma3 we additionally modify the architecture to include dropout after the gated product, in order to test sensitivity to stochastic masking.

For gated MLPs we train 3-layer models with hidden width 128 on MNIST and CIFAR-10, using Adam at 10^{-4} and no dropout. In addition to clean test accuracy, we evaluate adversarial robustness using the fast gradient sign method (FGSM), applying uniform input noise of 0.2 and 10 attack steps with a step size 0.015.

4.1 ViT-VAE ON CELEBA

We first examine a Vision Transformer VAE on CelebA. Figure 1 shows training losses with and without a post-gating bias (PGB). At early epochs the two variants behave similarly. However, loss spikes occur sooner and persist longer in the model without PGB. By contrast, PGB shortens the duration of spikes, leading to longer regions of stable descent and ultimately a better asymptotic outcome.

Although the effect is modest, it is consistent with our analysis: PGB provides a direct intercept that can absorb noisy fluctuations in the latent signal, allowing gradients for W_d to be driven by centered features. Given that the implementation requires only an in-place addition and negligible parameter cost, the observed improvement highlights a potentially useful stabilization mechanism for transformer-based VAEs.

4.2 LLM PRETRAINING

We next consider large language model pretraining. Here we test Phi3-mini and Gemma3 with four bias configurations: (1) no biases, (2) PGB only, (3) RMSNorm bias only, and (4) both biases. The Gemma3 architecture is also modified to include dropout after the gated product.

Figures 2 and 3 show that none of these variants makes a significant difference in training loss. This supports the hypothesis that the main benefit of PGB is not in attention-heavy autoregressive settings, but rather in architectures where stochastic noise interacts with latent encodings, as in VAEs.

4.3 STACKED GATED MLPs ON MNIST AND CIFAR-10

Finally, we evaluated stacked Gated MLPs (without attention) on MNIST and CIFAR-10 under adversarial attack. Here an unexpected result emerged: adding bias after gating consistently improved robustness to fast gradient sign method (FGSM) perturbations.

Figure 4 shows adversarial test accuracy on CIFAR-10 with three hidden layers of width 128. The model with PGB maintains a clear advantage in adversarial accuracy throughout training. On MNIST, the effect is also visible. With three hidden layers of width 256 (Figure 5), PGB raises

324
325
326
327
328
329
330
331
332
333
334
335
336
337
338
339
340
341
342
343
344
345
346
347
348
349
350
351
352
353
354
355
356
357
358
359
360
361
362
363
364
365
366
367
368
369
370
371
372
373
374
375
376
377

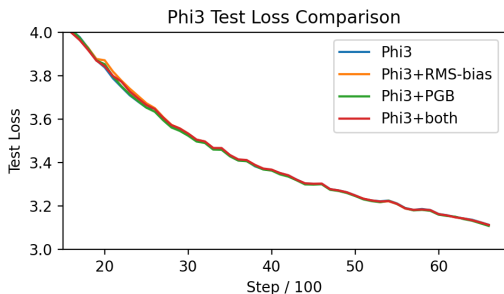


Figure 2: Test loss for Phi3 under four bias configurations. No significant effect is observed.

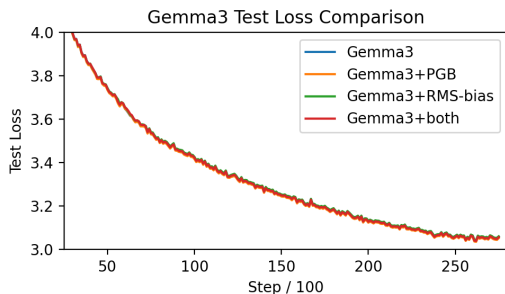


Figure 3: Test loss for Gemma3, including gated signal dropout, under four bias configurations. Again, no discernible effect.

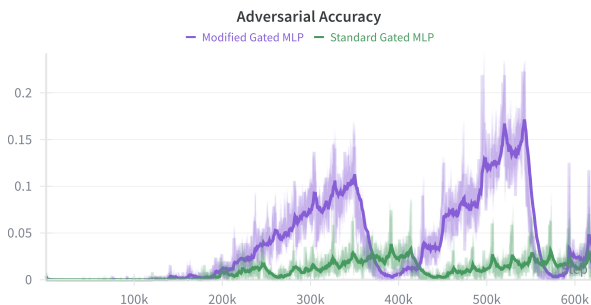


Figure 4: Adversarial accuracy during training of stacked Gated MLPs on CIFAR-10 (3 hidden layers, width 128). Models with PGB show stronger robustness to FGSM perturbations. All models are trained with Adam at learning rate 10^{-4} and no dropout.

adversarial accuracy from about 14% to 16%. Interestingly, in a smaller configuration with two hidden layers of width 64 (Figure 6), the benefit is much more pronounced, suggesting that PGB may provide stronger robustness in lower-capacity models.

Although preliminary, these results point to a potential role for explicit baselines in enhancing adversarial robustness in purely feed-forward gated networks. Further work is needed to test whether this effect generalizes beyond MNIST and CIFAR-10 and under stronger adversarial settings.

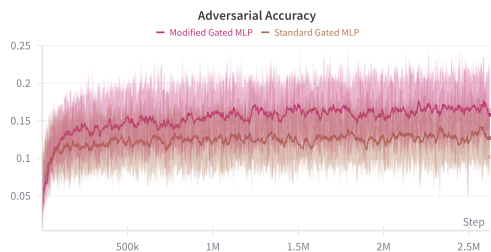


Figure 5: Adversarial accuracy on MNIST for stacked Gated MLPs with 3 hidden layers of width 256. Here PGB raises adversarial test accuracy from $\sim 14\%$ to $\sim 16\%$.

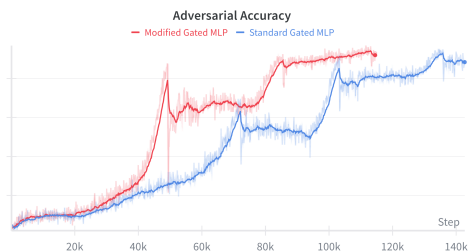


Figure 6: Adversarial accuracy on MNIST for stacked Gated MLPs with 2 hidden layers of width 64. PGB provides a pronounced robustness gain under FGSM perturbations compared to the bias-free baseline.

Grokking Phenomenon. Interestingly, the phenomenon of grokking, where neural networks gradually become more robust or generalize better long after achieving near-perfect training perfor-

mance, has been documented not only in algorithmic tasks but also in classification settings, including adversarial robustness (Humayun et al., 2024). In their formulation, delayed robustness arises as the network’s partition of linear regions migrates in input space, reducing sensitivity around training points and pushing non-linearities toward decision boundaries. From this perspective, PGB might act as an accelerant or enabler of this transition: by stabilizing the baseline routing early, it may speed up or strengthen the migration dynamics that underlie robust grokking. In our experiments, we observe that adversarial accuracy with PGB begins rising earlier and stabilizes higher than the bias-free counterpart, consistent with the view that PGB helps the network “grok robustness” sooner or more reliably.

5 DISCUSSION

Our study began from an empirical observation: introducing a post-gating bias (PGB) in a transformer VAE yielded a marked improvement in training stability and asymptotic loss. Motivated by this finding, we analyzed the role of PGB as restoring affine freedom after multiplicative gating. The analysis shows that without an explicit intercept, weight-only recovery of constant offsets is ill-conditioned, especially when dropout or latent noise perturb the mean direction of post-gating features. A trainable PGB both improves conditioning (by removing a rank-1 term from the Hessian) and reduces variance in gradients (by centering the masked regressor under dropout).

Our experiments provide preliminary but informative evidence:

- In a ViT-VAE on CelebA, PGB shortens loss spikes and yields a slight asymptotic benefit, consistent with the variance-reduction interpretation.
- In two LLM architectures (Gemma3 and Phi3-mini), adding PGB or other biases had no discernible effect, suggesting that autoregressive pretraining is less sensitive to this mechanism.
- In stacked Gated MLPs, PGB improved robustness to adversarial perturbations, an effect that deserves further study.

Taken together, these findings suggest that PGB is not tied solely to gated dropout, as we initially suspected, but may be most relevant in settings where stochastic noise interacts with latent encodings. The VAE experiments point in this direction, as PGB provides a baseline that absorbs fluctuations otherwise routed through ill-conditioned weights. We also observed that training instabilities were most visible early in optimization, when the learning rate is large, but largely vanish once the rate decays. This interplay between gradient variance, noise level, and learning-rate dynamics deserves closer study.

The analysis highlights that the main role of PGB is geometric: by removing a rank-1 outer product from the effective Hessian, it reduces the dominance of a single high-variance direction. This is less about expressivity (which is unchanged) and more about conditioning. From this perspective, PGB is a form of variance control: it reshapes the eigen-spectrum so that small perturbations are not amplified along unstable directions, making optimization smoother.

More broadly, our results caution against the common view that biases in transformers are entirely redundant. While this heuristic holds in attention-heavy LLMs, our analysis and experiments reveal that redundancy breaks down after multiplicative gating. Thus, what appears redundant at the level of expressivity may not be redundant at the level of optimization dynamics. Architectural simplifications motivated solely by redundancy arguments should therefore be revisited when conditioning or robustness are at stake.

Here, even a negligible-cost bias can alter the conditioning of the optimization problem, stabilizing features that would otherwise be transmitted through high-variance channels. In this sense, PGB is less about expressivity and more about *geometry*: it reshapes the spectrum of the effective Hessian in a way that lowers gradient variance and smooths descent.

The connection to robustness is particularly intriguing. Adversarial experiments on MNIST and CIFAR-10 show that PGB increases resistance to perturbations, echoing recent arguments that robustness eventually emerges in overparameterized models through continued training. By supplying

432 a baseline earlier in optimization, PGB may accelerate this “grokking” of robustness or improve its
433 asymptotic level, though confirming this will require systematic study.

434 Finally, we emphasize the practicality of the modification. In an era where the cost of architectural
435 changes is heavily scrutinized, PGB is effectively free: it can be implemented as an in-place addition,
436 introducing negligible parameter or memory overhead. This positions it as a candidate for inclusion
437 wherever training instabilities or robustness concerns are present. Future work should test PGB
438 more broadly: across multiple seeds and noise levels in VAEs, with spectral diagnostics that track
439 the evolution of Hessian eigenvalues, and in domains where latent noise is inherent, such as diffusion
440 autoencoders or scientific surrogate models.

441 In summary, PGB is not a universal fix, but it offers a small and elegant tool for mitigating the
442 interaction between gating, stochastic noise, and optimization stability. Its simplicity makes it easy
443 to adopt, and even modest improvements may prove valuable in practice.

444 6 CONCLUSION

445 We introduced Post-Gating Bias (PGB), a simple addition that restores affine freedom in transformer
446 MLPs after multiplicative gating. Our analysis shows that PGB improves conditioning and reduces
447 gradient variance under noise, without altering expressivity or incurring cost. Experiments confirm
448 that its impact is modest in autoregressive pretraining but beneficial in VAEs and gated MLPs,
449 where it stabilizes training and enhances robustness. These results highlight that even “redundant”
450 components can matter for optimization dynamics, suggesting that bias placement in gated networks
451 deserves renewed attention.

452 7 LIMITATIONS AND OUTLOOK

453 Our evaluation is limited in several respects. First, we have not presented the crystallographic VAE
454 case that originally motivated this study, as it depends on an unpublished architecture. Second, our
455 experimental trials are few in number, and further repetitions are needed to confirm robustness across
456 seeds, hyperparameters, and training regimes. The absence of gains in LLM pretraining underscores
457 the need to better delineate when and why PGB helps, and whether its role is confined to noisy latent
458 encoders or extends more broadly.

459 A further limitation concerns adversarial robustness and grokking. While our stacked gated MLP
460 experiments revealed improvements in adversarial accuracy, these results are preliminary and not yet
461 fully understood. We observed interactions with architectural depth and width, suggesting that the
462 benefits of PGB may depend on how baselines are represented across layers. Moreover, although
463 the robustness gains are suggestive of earlier grokking, consistent with recent work showing that
464 robustness eventually emerges with sufficient training, we cannot yet determine whether PGB ac-
465 celerates this process, improves asymptotic robustness, or simply interacts with optimization noise
466 in a coincidental way. Even in LLMs, the seemingly identical loss trajectories may conceal deeper
467 differences in the learned representations that our current evaluation does not capture.

468 That said, we believe this work is a useful contribution: it highlights a simple modification, grounded
469 in analysis, that can sometimes mitigate instability at effectively zero cost. We hope this encourages
470 further experimentation in diverse architectures, particularly in noisy or generative settings where
471 stable optimization remains a challenge.

472 **Use of LLMs.** We used ChatGPT (GPT-5, OpenAI) to assist with improving the clarity and consis-
473 tency of exposition, formatting LaTeX, and checking grammar. All ideas, experiments, and analyses
474 were conceived, designed, and executed by the authors.

475 REFERENCES

476 Marah Abdin, Joshua Ainslie, Chandramouli Anchuri, Anton Bakhtin, Shruti Bhosale, Ben Bo-
477 gin, Junteng Chen, Hao Cheng, Jonathan H Clark, et al. Phi-3 technical report. *arXiv preprint*
478 *arXiv:2404.14219*, 2024.

- 486 Huiwen Chang, Han Zhang, Lu Jiang, Ce Liu, William T. Freeman, Michael Rubinstein, and Kfir
487 Aberman. Maskgit: Masked generative image transformer. In *Proceedings of the IEEE/CVF*
488 *Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 11315–11325, 2022.
- 489 Aakanksha Chowdhery, Sharan Narang, Jacob Devlin, Maarten Bosma, Gaurav Mishra, Adam
490 Roberts, Paul Barham, Hyung Won Chung, Charles Sutton, Sebastian Gehrmann, et al. Palm:
491 Scaling language modeling with pathways. *arXiv preprint arXiv:2204.02311*, 2022.
- 492 Yann N Dauphin, Angela Fan, Michael Auli, and David Grangier. Language modeling with gated
493 convolutional networks. In *International conference on machine learning*, pp. 933–941. PMLR,
494 2017.
- 495 James Henderson and Jannis Fehr. A variational autoencoder for transformers with nonparametric
496 variational information bottleneck. *arXiv preprint arXiv:2207.13529*, 2022.
- 497 Zhiting Hu, Xuezhe Li, Percy Liang, and Eric P Xing. Recurrence boosts diversity! revisiting
498 recurrent latent variable in transformer-based vae for diverse text generation. *arXiv preprint*
499 *arXiv:2210.12409*, 2022.
- 500 Ahmed Imtiaz Humayun, Randall Balestrieri, and Richard Baraniuk. Deep networks always grok
501 and here is why. *arXiv preprint arXiv:2402.15555*, 2024.
- 502 Rafal Jozefowicz, Wojciech Zaremba, and Ilya Sutskever. An empirical exploration of recurrent
503 network architectures. In *International conference on machine learning*, pp. 2342–2350. PMLR,
504 2015.
- 505 Rafal Jozefowicz, Oriol Vinyals, Mike Schuster, Noam Shazeer, and Yonghui Wu. Exploring the
506 limits of language modeling. *arXiv preprint arXiv:1602.02410*, 2016.
- 507 Alex Krizhevsky. Learning multiple layers of features from tiny images. 2009. URL <https://api.semanticscholar.org/CorpusID:18268744>.
- 508 Minyoung Lee, Youngjune Kim, Sungjin Cho, and Wonyong Sung. Autoregressive image generation
509 using residual quantization. In *Proceedings of the AAAI Conference on Artificial Intelligence*,
510 volume 34, pp. 4350–4357, 2020.
- 511 A Ramesh and K Ramkumar. Mabvit: Modified attention block enhances vision transformers. *arXiv*
512 *preprint arXiv:2312.01324*, 2023.
- 513 Pौर्या Shamsolmoali, Masoumeh Zareapoor, et al. Vtae: Variational transformer autoencoder with
514 manifolds. *arXiv preprint arXiv:2304.00948*, 2023.
- 515 Noam Shazeer. Glu variants improve transformer. *arXiv preprint arXiv:2002.05202*, 2020.
- 516 David R So, Wojciech Mańke, Quoc V Le, Noam Shazeer, et al. Primer: Searching for efficient
517 transformers for language modeling. In *Advances in Neural Information Processing Systems*,
518 2021.
- 519 A Solera-Rico, R Murayama, et al. β -variational autoencoders and transformers for reduced-order
520 modelling of fluid flows. *arXiv preprint arXiv:2304.03571*, 2023.
- 521 Google DeepMind Team. Gemma 3: Open models for the open community. *arXiv preprint*
522 *arXiv:2412.19437*, 2024.
- 523 Hugo Touvron, Thibaut Lavril, Gautier Izacard, Herve Jegou, Damien Gravier, Pierric Boone, et al.
524 Llama: Open and efficient foundation language models. *arXiv preprint arXiv:2302.13971*, 2023.
- 525 Kai Wang, Xiaojun Wan, and Ming Yang. T-cvae: Transformer-based conditioned variational au-
526 toencoder. In *Proceedings of the 28th International Joint Conference on Artificial Intelligence*,
527 pp. 5273–5279, 2019.
- 528 Han Xiao, Kashif Rasul, and Roland Vollgraf. Fashion-mnist: a novel image dataset for benchmark-
529 ing machine learning algorithms, 2017.
- 530 Wei Zhang, Jun Li, et al. Graph-induced syntactic-semantic spaces in transformer-based variational
531 autoencoders. *arXiv preprint arXiv:2311.08579*, 2023.
- 532
533
534
535
536
537
538
539

APPENDIX A PROOFS FOR ANALYSIS SECTION

A.1 PROOF OF LEMMA 1 (MINIMAL-CHANGE EXACT-FIT)

We solve

$$\min_{\Delta} \|\Delta\|_F \quad \text{s.t.} \quad \Delta \mathbf{b}_{\text{in}} = \mathbf{b}_{\text{out}}.$$

Let $\mathbf{u} = \mathbf{b}_{\text{in}} / \|\mathbf{b}_{\text{in}}\|_2$. Any feasible Δ must satisfy $\Delta \mathbf{u} = \mathbf{b}_{\text{out}} / \|\mathbf{b}_{\text{in}}\|_2$. Among all linear maps that send \mathbf{u} to this value, the one with smallest Frobenius norm is the rank-1 map

$$\Delta^* = \frac{\mathbf{b}_{\text{out}} \mathbf{u}^\top}{\|\mathbf{b}_{\text{in}}\|_2} = \frac{\mathbf{b}_{\text{out}} \mathbf{b}_{\text{in}}^\top}{\|\mathbf{b}_{\text{in}}\|_2^2}.$$

Uniqueness follows from orthogonality of the Frobenius norm: any additional component off $\text{span}(\mathbf{b}_{\text{in}})$ increases $\|\Delta\|_F$ without improving feasibility. \square

A.2 PROOF OF REMARK (MINIMUM-MSE PERTURBATION)

We minimize

$$\min_{\Delta} \mathbb{E} \|(\mathbf{W}(\mathbf{x} + \mathbf{b}_{\text{in}}) + \mathbf{b}_{\text{out}}) - (\mathbf{W} + \Delta)(\mathbf{x} + \mathbf{b}_{\text{in}})\|_2^2,$$

where \mathbf{x} has mean zero and covariance Σ . Expanding,

$$\mathbb{E} \|\mathbf{W}(\mathbf{x} + \mathbf{b}_{\text{in}}) + \mathbf{b}_{\text{out}} - (\mathbf{W} + \Delta)(\mathbf{x} + \mathbf{b}_{\text{in}})\|_2^2 = \mathbb{E} \|(\mathbf{b}_{\text{out}} - \Delta \mathbf{b}_{\text{in}}) - \Delta \mathbf{x}\|_2^2.$$

Since $\mathbb{E}[\mathbf{x}] = 0$, this equals

$$\|\mathbf{b}_{\text{out}} - \Delta \mathbf{b}_{\text{in}}\|_2^2 + \text{tr}(\Delta \Sigma \Delta^\top).$$

This is a convex quadratic in Δ with unique minimizer given by normal equations

$$(\Sigma + \mathbf{b}_{\text{in}} \mathbf{b}_{\text{in}}^\top) \Delta^\top = \mathbf{b}_{\text{in}} \mathbf{b}_{\text{out}}^\top,$$

hence

$$\Delta^* = \mathbf{b}_{\text{out}} \mathbf{b}_{\text{in}}^\top (\Sigma + \mathbf{b}_{\text{in}} \mathbf{b}_{\text{in}}^\top)^{-1}.$$

Substituting back yields the minimal MSE formula in the main text. \square

A.3 PROOF OF PROPOSITION 1 (INTERCEPT REMOVAL IMPROVES CONDITIONING)

Let $\varphi \in \mathbb{R}^p$ denote post-gating features with mean μ and covariance Σ . The uncentered Hessian is

$$\mathbf{H}_{\text{no-bias}} = \mathbb{E}[\varphi \varphi^\top] = \Sigma + \mu \mu^\top,$$

while with explicit bias the Hessian reduces to $\mathbf{H}_{\text{pgb}} = \Sigma$.

By the Rayleigh quotient, for the normalized top eigenvector \mathbf{v} of Σ ,

$$\lambda_{\max}(\mathbf{H}_{\text{no-bias}}) \geq \mathbf{v}^\top (\Sigma + \mu \mu^\top) \mathbf{v} = \lambda_{\max}(\Sigma) + (\mathbf{v}^\top \mu)^2.$$

Thus $\lambda_{\max}(\mathbf{H}_{\text{no-bias}}) \geq \lambda_{\max}(\Sigma)$, with strict inequality when $\mu \notin \ker(\Sigma)$. Since both Hessians share the same nullspace, the condition number satisfies

$$\kappa(\mathbf{H}_{\text{no-bias}}) \geq \kappa(\Sigma),$$

with strict inequality unless $\mu \in \ker(\Sigma)$. \square

A.4 PROOF OF REMARK (CENTERING UPSTREAM REDUCES SECOND MOMENTS)

Suppose $\mathbf{X} \in \mathbb{R}^{n \times d}$ has tokenwise expectations $\bar{\mathbf{X}}$ (each row mean of \mathbf{X}), with global mean $\bar{\mu} = \frac{1}{n} \mathbf{1}^\top \bar{\mathbf{X}}$. Replacing \mathbf{X} by $\mathbf{X}' = \mathbf{X} - \mathbf{1} \bar{\mu}^\top$ changes the Frobenius norm:

$$\mathbb{E} \|\mathbf{X}'\|_F^2 = \mathbb{E} \|\mathbf{X}\|_F^2 - n \|\bar{\mu}\|_2^2.$$

Thus the expected squared norm is strictly reduced by $\|\bar{\mu}\|_2^2$ per token.

If $\mathbf{U} = \mathbf{X} \mathbf{W}_u$ and $\mathbf{G} = \mathbf{X} \mathbf{W}_g$, then

$$\mathbb{E} \|\mathbf{U}'\|_F^2 = \mathbb{E} \|\mathbf{U}\|_F^2 - \|\bar{\mu} \mathbf{W}_u\|_2^2, \quad \mathbb{E} \|\mathbf{G}'\|_F^2 = \mathbb{E} \|\mathbf{G}\|_F^2 - \|\bar{\mu} \mathbf{W}_g\|_2^2.$$

Since SiLU is 1-Lipschitz and monotone,

$$\mathbb{E} \|\mathbf{S}'\|_F^2 = \mathbb{E} \|\mathbf{U}' \odot \text{SiLU}(\mathbf{G}')\|_F^2 \leq \mathbb{E} \|\mathbf{U}'\|_F^2 \mathbb{E} \|\text{SiLU}(\mathbf{G}')\|_\infty^2,$$

which bounds the second moment of \mathbf{S} after centering. \square

A.5 PROOF OF LEMMA 2 (BASELINE MINIMIZES GRADIENT VARIANCE)

For row j , define

$$\mathbf{g}_i(j, u) = \mathbf{u}^\top \left(\frac{\partial \ell_i}{\partial \mathbf{W}_d} \right)_j^\top = (\mathbf{M}_{ij} \mathbf{S}_{ij} + \mathbf{b}_j) r_i, \quad r_i = \mathbf{u}^\top \boldsymbol{\gamma}_i.$$

By independence and $\mathbb{E}[\boldsymbol{\gamma}_i] = 0$,

$$\text{Var}(\mathbf{g}_i(j, u)) = \mathbb{E}[r_i^2] \text{Var}(\mathbf{M}_{ij} \mathbf{S}_{ij} + \mathbf{b}_j).$$

Now,

$$\mathbb{E}[\mathbf{M}_{ij} \mathbf{S}_{ij}] = p \boldsymbol{\mu}_j, \quad \text{Var}(\mathbf{M}_{ij} \mathbf{S}_{ij}) = p \sigma_j^2 + p(1-p) \boldsymbol{\mu}_j^2.$$

Hence

$$\text{Var}(\mathbf{M}_{ij} \mathbf{S}_{ij} + \mathbf{b}_j) = p \sigma_j^2 + (p \boldsymbol{\mu}_j + \mathbf{b}_j)^2,$$

and therefore

$$\text{Var}(\mathbf{g}_i(j, u)) = \mathbb{E}[r_i^2] (p \sigma_j^2 + (p \boldsymbol{\mu}_j + \mathbf{b}_j)^2).$$

This quadratic in \mathbf{b}_j is uniquely minimized at $\mathbf{b}_j^* = -p \boldsymbol{\mu}_j$, with minimum value $\mathbb{E}[r_i^2] p \sigma_j^2$. \square

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