# GVPO Explained: Novelties, Advantages, and Mathematical Differences

## Executive Summary

**GVPO (Group Variance Policy Optimization)** addresses the training instability issues of GRPO while providing stronger theoretical guarantees. The key innovation is incorporating the **analytical solution to KL-constrained reward maximization directly into gradient weights** through a clever **zero-sum weight constraint** that eliminates the intractable partition function.

---

## 1. Core Novelties of GVPO

### Novel 1: Zero-Sum Weight Constraint Eliminates Partition Function

**The Problem:** The optimal policy for KL-constrained reward maximization has a closed-form solution:

$$\pi^*(y|x) = \frac{1}{Z(x)}\pi_{\theta'}(y|x)e^{R(x,y)/\beta}$$

where $Z(x) = \sum_y \pi_{\theta'}(y|x)e^{R(x,y)/\beta}$ is computationally intractable (requires summing over all possible responses).

**GVPO's Solution:** By designing weights such that $\sum_{i=1}^{k} w_i = 0$, the partition function $\beta \log Z(x)$ becomes **invariant across responses and cancels out** in gradient computations:

$$\nabla_\theta L(\theta) = -\sum_{x,\{y_i\}} \sum_{i=1}^{k} w_i \nabla_\theta \log \frac{\pi_\theta(y_i|x)}{\pi_{\theta'}(y_i|x)} = -\sum_{x,\{y_i\}} \sum_{i=1}^{k} w_i \nabla_\theta \frac{R_\theta(x,y_i)}{\beta}$$

Since $\sum w_i = 0$, the $\beta \log Z(x)$ term disappears, making the method computationally tractable.

### Novel 2: Gradient Weights Based on Central Distance Differences

**GVPO's Weight Design:**

$$w_i = (R(x, y_i) - \bar{R}(x)) - \beta \left( \log \frac{\pi_\theta(y_i|x)}{\pi_{\theta'}(y_i|x)} - \overline{\log \frac{\pi_\theta}{\pi_{\theta'}}} \right)$$

where the bar notation denotes group average: $\bar{R}(x) = \frac{1}{k} \sum_{i=1}^{k} R(x, y_i)$.

**Physical Interpretation:** The weight is the **difference between actual reward central distance and implicit reward central distance**.

**Novel 3: Three Equivalent Loss Interpretations**

The paper elegantly shows GVPO's loss has three mathematically equivalent forms:

**(a) Negative Log-Likelihood View (Equation 9):**

$$\mathcal{L}_{\text{GVPO}}(\theta) = -\beta \sum_{x,\{y_i\}} \sum_{i=1}^{k} \left[ (R(x, y_i) - \bar{R}) - \beta \left( \log \frac{\pi_\theta(y_i|x)}{\pi_{\theta'}(y_i|x)} - \overline{\log \frac{\pi_\theta}{\pi_{\theta'}}} \right) \right] \log \pi_\theta(y_i|x)$$

**(b) Mean Squared Error View (Middle panel, Figure 1):**

$$\nabla_\theta \mathcal{L}_{\text{GVPO}} = \frac{1}{2} \nabla_\theta \sum_{x,\{y_i\}} \sum_{i=1}^{k} \left[ (R_\theta(x, y_i) - \bar{R}_\theta) - (R(x, y_i) - \bar{R}) \right]^2$$

**Key Insight:** Minimizing GVPO loss = minimizing **MSE between implicit and actual reward central distances**.

**(c) Reinforcement Learning View (Equation 14, =1):**

$$\nabla_\theta \hat{\mathcal{L}}_{\text{GVPO}} = -2\mathbb{E}_{x,y} \left[ (R(x, y) - \mathbb{E}_y R) \log \pi_\theta(y|x) + \text{Cov}(\log \pi_\theta, \log \pi_{\theta'}) - 0.5\text{Var}(\log \pi_\theta) \right]$$

Three components: 1. **Group-relative reward term**: Advantage maximization
2. **Covariance term**: Regularization preventing deviation from reference policy
3. **Variance term**: Entropy-like exploration encouragement

---

## 2. Mathematical Comparison: GVPO vs GRPO

**GRPO Loss (Equation 2):**

$$\mathcal{L}_{\text{GRPO}}(\theta) = - \sum_{x,y_1,\ldots,y_k} \sum_{i=1}^{k} \frac{R(x, y_i) - \text{Mean}(\{R(x, y_i)\})}{\text{Std}(\{R(x, y_i)\})} \log \pi_\theta(y_i|x)$$

**Key Differences:**

| Aspect | GRPO | GVPO |
|---|---|---|
| **Weight Formula** | $w_i = \frac{R(x,y_i) - \bar{R}}{\sigma_R}$ (standardized reward) | $w_i = (R(x, y_i) - \bar{R}) - \beta(\log \frac{\pi_\theta}{\pi_{\theta'}} - \overline{\log \frac{\pi_\theta}{\pi_{\theta'}}})$ |
| **Normalization** | Divides by standard deviation $\sigma_R$ | No std normalization (only centering) |

| Aspect | GRPO | GVPO |
|---|---|---|
| **Policy Dependency** | Weights independent of current policy | Weights depend on $\pi_\theta/\pi_{\theta'}$ ratio |
| **KL Constraint** | Applied externally (hyperparameter tuning) | **Built into gradient weights analytically** |
| **Zero-Sum Property** | Yes (due to centering) | Yes (by design) |

**Critical Mathematical Insight:**

GRPO's standardization **conflates prompt-level difficulty with reward signals** (cited in paper [17]). For example: - Hard prompt with rewards [8, 9, 10] → all responses get similar standardized scores - Easy prompt with rewards [1, 2, 9] → large standardized score differences

GVPO **removes std normalization** but adds the $\beta(\log \pi_\theta/\pi_{\theta'})$ term to directly encode the optimal policy structure.

---

## 3. Theoretical Advantages of GVPO

**Advantage 1: Unique Optimal Solution (Theorem 3.1)**

**GVPO Guarantee:**

$$\text{argmin}_\theta \hat{\mathcal{L}}_{\text{GVPO}}(\theta) = \pi^*(y|x) = \frac{1}{Z(x)}\pi_{\theta'}(y|x)e^{R(x,y)/\beta}$$

**Uniqueness** is proven by showing: 1. When $\pi_\theta = \pi^*$, the loss equals 0 (minimum achieved) 2. Any other policy yields loss > 0 (contradiction proof in Appendix B.1)

**Why This Matters:** - **DPO fails this**: Due to Bradley-Terry model limitations [3, 11], DPO may converge to suboptimal policies - **GRPO lacks this**: No theoretical guarantee of convergence to KL-constrained optimum

**Advantage 2: Flexible Sampling Distributions (Corollary 3.2)**

**GVPO's Condition:** Theorem 3.1 holds for **any sampling distribution $\pi_s$** satisfying:

$$\forall x, \{y|\pi_{\theta'}(y|x) > 0\} \subseteq \{y|\pi_s(y|x) > 0\}$$

**Translation:** As long as $\pi_s$ covers all responses that the reference policy could generate, GVPO maintains theoretical guarantees.

**Comparison with GRPO/PPO:**

| Method | Sampling Requirement | Problem |
|--------|---------------------|---------|
| **PPO** | On-policy $(\pi_s = \pi_\theta)$ | Low sample efficiency, requires fresh trajectories |
| **GRPO** | Uses importance sampling $\frac{\pi_\theta}{\pi_{\theta_{\text{old}}}}$ | Gradient explosion when policies diverge; requires clipping (introduces bias) |
| **GVPO** | Any $\pi_s$ satisfying mild condition | **No importance sampling, no explosion risk** |

**Mathematical Detail:** Policy gradient methods require:

$$\nabla_\theta[\mathbb{E}_{x,y\sim\pi_\theta}[R(x,y)]-\text{DKL}[\pi_\theta||\pi_{\theta_{\text{old}}}]] = \mathbb{E}_{x,y\sim\pi_\theta}\left[\left(R - \log\frac{\pi_\theta}{\pi_{\theta_{\text{old}}}} - 1\right)\nabla_\theta\log\pi_\theta\right]$$

Off-policy estimation uses importance sampling (Equation 16):

$$\mathbb{E}_{x,y\sim\pi_{\theta_{\text{old}}}}\left[\frac{\pi_\theta(y|x)}{\pi_{\theta_{\text{old}}}(y|x)}\left(R - \log\frac{\pi_\theta}{\pi_{\theta_{\text{old}}}} - 1\right)\nabla_\theta\log\pi_\theta\right]$$

The ratio $\frac{\pi_\theta}{\pi_{\theta_{\text{old}}}}$ can explode $\rightarrow$ gradient clipping needed.

**GVPO's Advantage:** By using the zero-sum property and central distances, GVPO's gradient becomes:

$$\mathbb{E}_{x,y\sim\pi_s}\left[\left(R - \log\frac{\pi_\theta}{\pi_{\theta'}} - \mathbb{E}_{y\sim\pi_s}\left(R - \log\frac{\pi_\theta}{\pi_{\theta'}}\right)\right)\nabla_\theta\log\pi_\theta\right]$$

**No importance sampling ratio** appears in the gradient!

**Advantage 3: Unbiased and Consistent Estimator (Theorem 3.4)**

The empirical loss with finite samples is:

$$\frac{1}{|D|}\sum_{(x,\{y_i\})\in D}\frac{1}{k-1}\sum_{i=1}^{k}\left[(R_\theta(x,y_i) - \bar{R}_\theta) - (R(x,y_i) - \bar{R})\right]^2$$

**Note the** $\frac{1}{k-1}$ **factor** (not $\frac{1}{k}$) — this is the **Bessel correction** for unbiased variance estimation.

**Why This Matters:** - With small $k$ (few samples per prompt), bias becomes significant - Corollary 3.5 extends this to **variable** $k(x)$ **per prompt**, enabling mixed-source datasets

---

## 4. Algorithm Comparison

**Algorithm 1 (GVPO) vs GRPO:**

```
GVPO:
1. Sample k responses {yi} ~ s(·|x)
2. Compute weights: wi = (R(x,yi) - R) - (log(/') - log(/'))
3. Update: minimize - wi log (yi|x)


GRPO:
1. Sample k responses {yi} ~ _old(·|x)
2. Compute weights: wi = (R(x,yi) - R) / R
3. Update: minimize - wi log (yi|x)
4. Apply gradient clipping + KL penalty
```

**Key Implementation Difference (Listing 1):**

GVPO only changes GRPO's loss computation by:

```
# GRPO:
advs = (R - R.mean()) / R.std()   # Standardization
loss = -scores * advs


# GVPO:
advs = (R - R.mean()) - beta * ((scores_new - scores_new.mean())
                              - (scores_old - scores_old.mean()))
loss = -beta * scores * advs / (k-1)  # Note: k-1 for unbiased estimator
```

---

## 5. Empirical Performance (Table 1)

| Model | AIME2024 | AMC | MATH500 | Minerva | OlympiadBench |
|---|---|---|---|---|---|
| Base (Qwen2.5-Math-7B) | 14.68 | 38.55 | 64.00 | 27.20 | 30.66 |
| +GRPO | 14.79 | 55.42 | **80.00** | 41.17 | 42.07 |
| +Dr.GRPO | 16.56 | 48.19 | 81.20 | 44.48 | 43.40 |

| Model | AIME2024 | AMC | MATH500 | Minerva | OlympiadBench |
|-------|----------|-----|---------|---------|---------------|
| **+GVPO** | **20.72** | **62.65** | **83.80** | **45.95** | **46.96** |

**Observations:** - GVPO achieves **best performance across all 5 benchmarks** - **40% relative improvement** on AIME2024 over GRPO (14.79 → 20.72) - Particularly strong on complex reasoning tasks (AIME, OlympiadBench)

**Ablation Study Insights:**

**Figure 2 ( sensitivity):** - GVPO shows **little performance fluctuation** across  [0.01, 0.5] - Suggests **robustness to hyperparameter tuning** (unlike GRPO's high sensitivity)

**Figure 3 (Scaling with k):** - GVPO **consistently outperforms GRPO** for all k  [2, 32] - **Superior scalability**: GVPO on 1.5B model with k=32 matches 7B model performance - **Inference cost reduction**: Can use smaller models with more samples

**Figure 4 (Off-policy sampling s):** - Tests mixing historical responses with current policy samples - GVPO maintains **robust performance** with ratios from 0:8 to 4:4 (historical:current) - Validates Corollary 3.2's theoretical guarantee

---

## 6. Limitations and Future Work

**Acknowledged Limitations:**

1. **Computational cost**: Still requires sampling k responses per prompt
2. **Reward model quality**: Performance depends on accurate R(x,y)
3. **Hyperparameter** : Though robust, still requires selection

**Unexplored Connections:**

- Integration with exploration strategies from classical RL
- Extension to continuous action spaces
- Multi-modal reward signals

---

## 7. Summary: Why GVPO is Better

| Criterion | GRPO | GVPO |
|-----------|------|------|
| **Training Stability** | Documented instability [34, 16] | Implicit regularization via Cov/Var terms |

| Criterion | GRPO | GVPO |
|---|---|---|
| **Hyperparameter Sensitivity** | High (clip threshold, KL coeff) | Robust to variations |
| **Theoretical Guarantee** | No convergence to optimal policy | Unique optimal = KL-constrained optimum |
| **Sampling Flexibility** | Uses importance sampling | Any s (no IS needed) |
| **Normalization Bias** | Std normalization conflates difficulty | Only centering (no std division) |
| **Gradient Explosion** | Requires clipping | No IS ratio $\rightarrow$ inherently stable |
| **Performance** | Baseline | **Best across all benchmarks** |

## Bottom Line

GVPO's core innovation is **operationalizing the closed-form optimal policy** through a mathematically elegant **zero-sum weight design** that: 1. Eliminates the intractable partition function 2. Embeds KL constraints directly into gradients 3. Enables off-policy training without importance sampling 4. Guarantees convergence to the unique optimal policy

This makes GVPO a **theoretically principled AND empirically superior** alternative to GRPO for LLM post-training.