

# 面向特定任务的大语言模型离散提示词进化方法 [你 的名字] [你的学院/单位名称]

**摘要** 大语言模型 (LLM) 在特定领域任务 (如数学推理、代码生成) 中的表现高度依赖于精心设计的 System Prompt。然而, 人工编写“最佳提示词”面临搜索空间爆炸 (Vocab Size > 50k) 与梯度不可导的双重挑战。本文提出一种基于语义子空间与 GRPO 的自动提示工程框架。该方法摒弃了全词表搜索的低效路径, 首先通过构建领域指令子集 (Sub-Vocab) 来锁定“专家语义空间”; 随后, 利用低维索引映射机制将优化参数量压缩至 KB 级别; 最后, 采用群组相对策略优化 (GRPO) 算法, 在无需训练价值网络 (Critic) 的前提下, 由任务数据集驱动模型自动“进化”出最优离散指令组合。该方法实现了低成本、高可解释性的特定任务领域适配。

**关键词** 特定任务适配; GRPO; 语义子空间; 离散优化; 查表映射

## 1 问题背景与核心直觉

### 1.1 不仅是“找咒语”, 而是“领域适配”

通用的 Prompt (如“你是一个有用的助手”) 无法激发模型在特定专业领域的潜力。我们的目标不是寻找一句通用的“神奇咒语”, 而是针对给定的某一类具体任务 (Task  $T$ ), 自动找到该任务下的最优指令  $P^*$ 。例如, 针对数学任务, 系统应自动发现类似“逐步拆解”、“反向验算”等关键指令的组合。

### 1.2 核心直觉: 受限空间内的概率进化

LLM 的 Prompt 搜索像是在 5 万个字的字典里大海捞针。我们的核心发现是: 并不需要尝试所有的词, 只需要在少数“关键指令词”中找到正确的组合。我们的方案遵循“缩小赌盘 -> 优胜劣汰”的逻辑:

1. **缩小赌盘:** 我们构建一个仅包含数百个高优指令词的“精英子词表”, 排除无关词汇的干扰。
2. **按单抓药:** 系统输出不再是具体的词, 而是子词表中的索引编号, 通过查表还原为文本。
3. **优胜劣汰:** 利用强化学习 (GRPO), 让那些能提高任务分数的指令组合“存活”下来。

## 2 方法设计

### 2.1 1. 数据结构: 静态映射与动态参数

**A. 领域子词表与映射表 ( $\mathcal{T}_{\text{map}}$ ):** 针对目标任务 (如编程), 筛选  $N = 300$  个高频指令词 (code, debug, function...), 构建静态映射表:

$$\mathcal{T}_{\text{map}} = [\text{id}_{\text{code}}, \text{id}_{\text{debug}}, \dots, \text{id}_{\text{optimize}}]$$

这相当于为该任务圈定了一个“专家词汇库”。

**B. 参数化生成矩阵 ( $\theta$ ):** 设定 Prompt 长度为  $L$  (如 20)。我们训练一个极小的矩阵  $\theta \in \mathbb{R}^{L \times N}$ 。其中  $\theta_{l,i}$  表示在第  $l$  个位置选择第  $i$  个指令词的非归一化概率。优势: 相比全词表方案 ( $L \times 50000$ ), 参数量减少了 99.4%, 仅需训练约 6000 个参数。

## 2.2 2. 前向生成过程 (Forward Process)

对于 Prompt 的每一个位置  $l \in \{1, \dots, L\}$ , 生成过程如下:

1. 概率计算: 对参数矩阵的第  $l$  行进行 Softmax:

$$\pi_l(z) = \text{Softmax}(\theta_l), \quad z \in \{0, \dots, N - 1\}$$

2. 离散采样: 根据分布  $\pi_l$  采样得到索引  $z_l$ 。

3. 查表映射: 将索引转换为真实 Token:  $t_l = \mathcal{T}_{\text{map}}[z_l]$ 。

4. 序列组装: 得到最终文本 Prompt  $P = [t_1, \dots, t_L]$ 。

## 2.3 3. 优化算法: GRPO (Group Relative Policy Optimization)

由于离散采样不可导, 我们采用 DeepSeek-R1 同款的 GRPO 算法, 无需 Critic 网络, 直接利用组内相对优势进行更新。

---

### Algorithm 1 特定任务提示词进化算法 (Subspace-GRPO)

---

**Require:** 任务数据集  $\mathcal{D}$ , 组采样数  $G = 16$ , 学习率  $\alpha$

```
1: 初始化参数  $\theta \leftarrow \mathcal{N}(0, 0.01)$ 
2: while not converged do
3:   从  $\mathcal{D}$  中抽取一批任务输入  $X$ 
4:   Step 1: 组采样 (Rollout)
5:   基于  $\theta$  采样  $G$  个不同的索引序列  $\{Z^{(1)}, \dots, Z^{(G)}\}$ 
6:   查表还原为 Prompt  $\{P^{(1)}, \dots, P^{(G)}\}$ 
7:   Step 2: 任务评估 (Evaluation)
8:   for  $g = 1$  to  $G$  do
9:     将  $P^{(g)} + X$  输入 LLM, 获得输出
10:    计算任务奖励  $r_g$  {如: 代码通过率、答案准确性}
11:   end for
12:   Step 3: 优势计算 (Advantage)
13:   计算组内平均分  $\mu_r$  和标准差  $\sigma_r$ 
14:    $A_g = (r_g - \mu_r) / (\sigma_r + \epsilon)$ 
15:   Step 4: 策略更新 (Update)
16:   最大化目标函数:  $J(\theta) = \frac{1}{G} \sum_g \left( \sum_l \log \pi_{\theta}(z_l^{(g)}) \cdot A_g \right)$ 
17:    $\theta \leftarrow \theta + \alpha \nabla J(\theta)$ 
18: end while
19: return 最优 Prompt  $P^*$ 
```

---

## 3 预期价值与应用

1. 自动化领域适配: 不需要人工专家, 系统自动根据 GSM8K 数据集进化出数学专用指令, 根据 HumanEval 进化出代码专用指令。

2. **完全可解释**: 最终产出的是人类可读的文本 (如”Think step-by-step”), 而非不可解释的连续向量, 便于分析模型行为。
3. **极低资源消耗**: 作为一个独立外挂模块, 仅需维护数 KB 的参数, 单张消费级显卡即可完成训练。

## 4 结论

本文提出了一种融合了“先验知识约束”与“强化学习探索”的提示词优化新范式。通过 Task-Specific Sub-Vocab 锁定语义范围, 利用 Mapping 机制降低计算维度, 结合 GRPO 实现高效进化, 我们为 LLM 的领域落地提供了一种低成本、高可用的自动适配工具。