

# PR #41775 完整报告

vllm-project/vllm

[Model Runner V2] FP32 gumbel sampling.

合并时间: 2026-05-16 00:20

原文链接: <http://prhub.com.cn/vllm-project/vllm/pull/41775>

## 执行摘要

- 一句话: Gumbel 采样默认使用 FP32 以提升性能
- 推荐动作: 值得精读学习如何在 Triton 内核中安全切换 FP32/FP64 并处理边界值; 以及从环境变量演化到引擎标志的设计决策过程, 体现了代码的健壮性和可维护性。

## 功能与动机

根据 PR body, H20 上的 profile 结果促使优化 Gumbel 采样精度行为。FP64 的使用会降低采样吞吐, 而经验上 FP32 对于 Gumbel-max 足够。目标是使 FP64 可选, FP32 成为默认快速选项, 并避免 Uniform 为零导致的数值问题。

## 实现拆解

1. 添加配置项: 在 vllm/config/model.py 的 ModelConfig 中新增 use\_fp64\_gumbel 布尔字段, 默认关闭, 不影响计算图。
2. 暴露 CLI 参数: 在 vllm/engine/arg\_utils.py 的 EngineArgs 中添加对应字段, 并注册 --use-fp64-gumbel CLI 参数, 将其传入 ModelConfig。
3. 修改 Gumbel 采样内核: 在 vllm/v1/worker/gpu/sample/gumbel.py 中新增 \_FP32\_TINY clamp 常量, 为 gumbel\_block\_argmax 和 \_gumbel\_sample\_kernel 添加 USE\_FP64 编译时常量参数, 根据其值选择 FP64/FP32 路径, 并调整 local\_max 张量的 dtype。
4. 修改拒绝采样内核: 在 vllm/v1/worker/gpu/spec\_decode/rejection\_sampler\_utils.py 中为 \_resample\_kernel 添加 USE\_FP64 参数, 并传递给 gumbel\_block\_argmax, 同时调整 resampled\_local\_max 的 dtype。
5. 传播配置: 在 Sampler、ModelRunner 以及推测解码模块的 RejectionSampler 和 EagleSpeculator 中接收并传递 use\_fp64\_gumbel, 确保内核能获得该配置。

关键文件:

- vllm/v1/worker/gpu/sample/gumbel.py (模块 采样内核; 类别 source; 类型 dependency-wiring; 符号 \_FP32\_TINY, gumbel\_block\_argmax, \_gumbel\_sample\_kernel, gumbel\_sample): 核心采样内核, 实现了 FP32/FP64 条件分支和 clamp 逻辑
- vllm/v1/worker/gpu/spec\_decode/rejection\_sampler\_utils.py (模块 推测解码; 类别 source; 类型 core-logic; 符号 \_resample\_kernel, rejection\_sample): 推测解码中的拒绝采样内核, 传递 USE\_FP64 参数并调整输出 dtype

- vllm/config/model.py (模块 模型配置; 类别 source; 类型 data-contract; 符号 ModelConfig.use\_fp64\_gumbel) : 新增 use\_fp64\_gumbel 配置字段, 作为单一数据源
- vllm/engine/arg\_utils.py (模块 引擎参数; 类别 source; 类型 core-logic; 符号 EngineArgs.use\_fp64\_gumbel) : 添加 --use-fp64-gumbel CLI 参数并传递到 ModelConfig
- vllm/v1/worker/gpu/sample/sampler.py (模块 采样器; 类别 source; 类型 core-logic; 符号 Sampler.init) : Sampler 类接受 use\_fp64\_gumbel 并传递给 gumbel\_sample
- vllm/v1/worker/gpu/model\_runner.py (模块 模型运行器; 类别 source; 类型 data-contract; 符号 ModelRunner.init) : ModelRunner 初始化时传递 use\_fp64\_gumbel 到 Sampler
- vllm/v1/worker/gpu/spec\_decode/eagle/speculator.py (模块 推测解码器; 类别 source; 类型 core-logic) : EagleSpeculator 传递 use\_fp64\_gumbel 到 rejection\_sample
- vllm/v1/worker/gpu/spec\_decode/rejection\_sampler.py (模块 拒绝采样器; 类别 source ; 类型 core-logic) : RejectionSampler 接收并传递 use\_fp64\_gumbel

关键符号: gumbel\_block\_argmax, \_gumbel\_sample\_kernel, gumbel\_sample, \_resample\_kernel, rejection\_sample, Sampler.init, ModelRunner.init, create\_model\_config

## 关键源码片段

### vllm/v1/worker/gpu/sample/gumbel.py

核心采样内核, 实现了 FP32/FP64 条件分支和 clamp 逻辑

```
# SPDX-License-Identifier: Apache-2.0
# SPDX-FileCopyrightText: Copyright contributors to the vLLM project
import torch
from vllm.triton_utils import HAS_TRITON, tl, triton

# 最小的正规 FP32 正值, 用于 clamp 均匀分布采样值,
# 避免 `log(u)` 产生 -inf, 从而保证 `-log(-log(u))` 始终有限。
# Triton 要求 `@triton.jit` 中访问的全局变量需用 `tl.constexpr(...)` 包裹,
# 但仅在 Triton 可用时才可调 (CPU worker 路径下 `tl` 为占位符, 调用会崩溃)。
_FP32_TINY = (
    tl.constexpr(float.fromhex("0x1p-126")) if HAS_TRITON else float.fromhex("0x1p-126")
)

@triton.jit
def gumbel_block_argmax(
    logits,
    block,
    mask,
    token_idx,
    expanded_idx_mapping_ptr,
    temp_ptr,
    seeds_ptr,
```

```

pos_ptr,
processed_logits_ptr,
processed_logits_stride,
processed_logits_col_ptr,
vocab_size,
APPLY_TEMPERATURE: tl.constexpr,
USE_FP64: tl.constexpr, # 新增参数: 是否使用 FP64 (默认 FP32)
):
req_state_idx = tl.load(expanded_idx_mapping_ptr + token_idx)
temp = tl.load(temp_ptr + req_state_idx).to(tl.float32)
if temp != 0.0 and APPLY_TEMPERATURE:
    logits = logits / temp

# ... 处理 processed_logits 存储 (省略) ...

# FP32 是默认规约精度; FP64 在 H100/Ada/Blackwell 上吞吐低 32-64 倍,
# 且经验上对于 Gumbel-max 不可区分。
if USE_FP64:
    logits = logits.to(tl.float64)
if temp != 0.0:
    seed = tl.load(seeds_ptr + req_state_idx)
    pos = tl.load(pos_ptr + token_idx)
    gumbel_seed = tl.randint(seed, pos)

    if USE_FP64:
        # 使用 64-bit 随机数构建 FP64 uniform
        u = tl_rand64(gumbel_seed, block, includes_zero=False)
    else:
        # 使用 FP32 uniform, 并 clamp 防止零
        u = tl.rand(gumbel_seed, block)
        u = tl.maximum(u, _FP32_TINY)
    gumbel_noise = -tl.log(-tl.log(u))

    logits = tl.where(mask, logits + gumbel_noise, float("-inf"))

value, idx = tl.max(logits, axis=0, return_indices=True)
return value, idx

```

## 评论区精华

- WoosukKwon建议将临时环境变量改为正式的引擎标志 `--use-fp64-gumbel`, 作者采纳。
- gemini-code-assist[bot]指出 `_gumbel_sample_kernel` 未添加 `USE_FP64` 参数会导致运行时 `TypeError`, 该问题在后续修复中解决。
- TheEpicDolphin要求确认 `draft acceptance rates` 在使用 FP32 后不会退化, 作者确认无问题。同时提出了一个代码风格 nit (将 `gumbel_noise` 计算提取出分支), 作者已修改。
- TheEpicDolphin还注意到一个多余文件 (pre-commit 输出) 被误加入, 作者移除。
- 将环境变量改为引擎标志 (design): 作者同意并修改为引擎标志。

- 缺少 USE\_FP64 参数导致 runtime error (correctness): 作者在后续修复中添加了该参数。
- 确认 FP32 对 draft acceptance 无退化 (testing): 作者确认无问题, 审查者表示信任。
- gumbel\_noise 计算可以提出分支 (style): 作者修改, 将计算移出分支。
- 误加多余文件 (other): 作者移除了该文件。

## 风险与影响

- 风险:
  - FP32 数值风险: 较低的随机精度可能在极少数情况下导致 Gumbel 采样出现 tie, 但通过 `_FP32_TINY clamp` 确保数值稳定, 且经验上差异不可感知。用户可通过 `--use-fp64-gumbel` 回退。
  - 默认行为变更: 默认使用 FP32 可能影响依赖原有 FP64 结果的用户, 提供了显式启用 FP64 的选项来兼容。
  - 测试缺口: 本 PR 未包含单元测试, 仅依赖手动验证和 profile。
  - 参数传递遗漏: 涉及多个内核文件, 可能引入参数传递遗漏, 但通过 Review 已捕获。
- 影响:
  - 用户影响: 默认 FP32 将提升采样阶段速度, 降低延迟, 对推理吞吐有益。需要使用精确 FP64 的用户可通过 `--use-fp64-gumbel` 启用。
  - 系统影响: 影响 V2 model runner 的采样和推测解码路径, 其他模块不受影响。
  - 团队影响: 该 PR 提供了可配置精度选项, 便于未来平衡性能与精度。
  - 风险标记: 核心路径变更, FP32 精度风险, 缺少测试覆盖

## 关联脉络

- 暂无明显关联 PR