

К использованию LLM в системах жёсткого реального времени

Старолетов Сергей Михайлович, кандидат физико-математических наук, доцент кафедры прикладной математики АлтГТУ, электронная почта: serg_soft@mail.ru

Широкое внедрение искусственного интеллекта для оценки текущего состояния переменных системы в промышленные робототехнические устройства и авионику сдерживается отсутствием формальных гарантий времени выполнения алгоритмов, ведь нормативные документы для доверенного программного обеспечения устройств (DO-178C, ISO 26262) требуют доказанного ограничения worst-case execution time. Предлагаемая работа – шаг к тому, чтобы показать, что современные облегченные языковые модели могут быть адаптированы к парадигме жёсткого реального времени.

В серии работ автор исследует свободную микроядерную ОС РОК [1], реализующую принцип разделов по стандарту ARINC 653, как платформу для экспериментов, поскольку она обеспечивает концепт строгой временной и пространственной изоляции компонентов – необходимое условие для будущей сертификации. В качестве LLM модели взята квантованная версия TinyStories (15 миллионов параметров), а в качестве кодовой базы LLM – код проекта llama2.c [3], способный работать на слабом железе IoT вроде одноплатных компьютеров.

Цель работы – создать демо, моделирующего работу «умного советчика» по параметрам окружающей среды. Каждое заданное количество миллисекунд система должна:

- Считать показания четырёх датчиков или их эмуляторов (температура, влажность, давление, скорость ветра).
- Преобразовать их в лингвистические оценки (вроде «жарко», «ветрено», «сухо»).
- Сформулировать запрос на естественном языке и передать его LLM внутри системы.
- Сгенерировать текстовую рекомендацию не позднее установленного дедлайна, даже если модель не успела выдать самый лучший ответ.

Модель TinyStories (dim=288, hidden_dim=768, n_layers=6, n_heads=6, vocab_size=32000), была квантована до int8 (для уменьшения памяти и поддержки вычислений в целых числах), преобразована программой в байты данных внутри .h файла и встроена получившимся статическим массивом в код elf-файла раздела (размер ~15 МБ). Прямой проход инференции состоит из последовательности трансформерных блоков. Время обработки одного слоя постоянно. Это даёт естественный квант для контроля времени – после завершения каждого слоя вызывается `rok_time_get()`, и, если текущее время превысило дедлайн, генерация завершается. Такой подход кооперативной проверки на границах слоёв гарантирует, что поток добровольно уступит управление не позже чем через время обработки одного слоя LLM.

Для сборки и успешной работы C-кода проекта llama2.c в среде РОК необходимо было сделать несколько решений портирования. Оригинальный код после загрузки квантованных весов сразу де-квантизует всю таблицу эмбедингов токенов, требуя для TinyStories около 36 дополнительных МБ (32000 * 288 * 4 байта). Под управлением РОК такой единовременный запрос к менеджеру кучи приводил к отказу. Поэтому решено было полностью отказаться от хранения де-квантизированной таблицы и вычислять эмбединг каждого токена «на лету» при каждом прямом проходе. В целом, можно вообще отказаться от динамического выделения памяти (это в коде делается в основном до начала работы при инициализации LLM), заменив

все это на статические массивы и арифметику указателей. Далее, хотя и в РОК есть libm для работы с математикой и часть libc, но было недоступны некоторые функции вроде qsort, bsearch, sprintf и так далее. Они были добавлены или код был переделан.

В итоге, в решении свободно доступном на GitHub [4], внутри одного раздела РОК организованы три потока, синхронизированные с помощью семафоров: (1) sensor_job – периодически опрашивает (для демо эмулирует с помощью случайной генерации и учета предыдущих значений по методу скользящего среднего) датчики температуры, скорости ветра, давления и влажности; (2) estimator_job – преобразует числовые показатели в строковые категории; (3) llm_job – формирует промпт, запускает инференцию (вывод) и выводит ответ. Результат можно посмотреть на Рис. 1.

```
[Estimator] classified temp=26 pressure=735 humidity=53 wind_speed=10
--- scheduling partition: 0, low:0, high:10
--- Scheduling processor: 0
    scheduling thread 2 (priority 42)
elected 2 !!!
    non-ready: 0 (1/stopped), 1 (40/waiting), 3 (42/lock), 4 (0/stopped), 5 (0/stopped), 6 (0/stopped), 7 (0/stopped), 8 (0/stopped)
[Sensor] get temperature 25 C
[Sensor] get humidity 53 %
[Sensor] get pressure 733 mm
[Sensor] get wind speed 8 m/s
--- scheduling partition: 0, low:0, high:10
--- Scheduling processor: 0
    scheduling idle thread
    non-ready: 0 (1/stopped), 1 (40/waiting), 2 (42/waiting), 3 (42/lock), 4 (0/stopped), 5 (0/stopped), 6 (0/stopped), 7 (0/stopped)
--- scheduling partition: 0, low:0, high:10
--- Scheduling processor: 0
    scheduling thread 2 (priority 42)
elected 2 !!!
    non-ready: 0 (1/stopped), 1 (40/waiting), 3 (42/lock), 4 (0/stopped), 5 (0/stopped), 6 (0/stopped), 7 (0/stopped), 8 (0/stopped)
--- scheduling partition: 0, low:0, high:10
--- Scheduling processor: 0
    scheduling idle thread
    non-ready: 0 (1/stopped), 1 (40/waiting), 2 (42/lock), 3 (42/lock), 4 (0/stopped), 5 (0/stopped), 6 (0/stopped), 7 (0/stopped)
--- scheduling partition: 0, low:0, high:10
--- Scheduling processor: 0
    scheduling thread 1 (priority 40)
elected 1 !!!
    non-ready: 0 (1/stopped), 2 (42/lock), 3 (42/lock), 4 (0/stopped), 5 (0/stopped), 6 (0/stopped), 7 (0/stopped), 8 (0/stopped)
[LLM] generation:
Temperature is hot and pressure is comfortable and wind is windy and humidity is comfortable. What to do?
Her mom had a surprise for her. She brought out a bag of colorful popcorn and said, "Let's all eat this popcorn together!"
Alice was very excited and said, "Yay!" They all shared the popcorn and enjoyed it in the sunshine.
[.]
[TRUNCATED at 91 tokens due to deadline]
[LLM] time for generate: 83422731
```

Рис.1 – Работа РОК LLM демо [4] (модель Stories 15M_Q8) на qemu-system-i386

В качестве продолжения возможна реализация полнофункционального демо с двумя разделами (партициями) – одним для LLM, а другим для датчиков и запуск его на устройстве. Модель уже была проверена на микрокомпьютере семейства Raspberry Pi и показала достаточно слабую скорость генерации – 2.18 токена/секунду, но это и нормально для чистого C-кода без оптимизаций. Но и получается, что LLM модель должна работать в секундных дедлайнах, выполняя долгосрочные прогнозы.

Также, простое завершение генерации по времени может оставить ответ синтаксически незавершённым. Можно реализовать более элегантный двухступенчатый механизм. За заданное в процентах до истечения дедлайна время может активироваться мягкий дедлайн. Модель при этом продолжает генерировать токены, но ей будет разрешено выдать не более заданного небольшого количества дополнительных токенов. Если в течение этого окна встречается знак окончания предложения – генерация немедленно останавливается и ответ получается грамматически законченным, иначе продолжается до истечения дедлайна и ответ получается как есть.

В итоге, можно констатировать, что задача может быть решена со свойственными таким моделям ограничениями и нужно искать случаи ее полезного использования.

Литература

1. POK kernel. Safe and secure for safety-critical systems. URL: <https://github.com/pok-kernel/pok>
2. Tiny Stories. URL: <https://huggingface.co/datasets/roneneldan/TinyStories> , <https://arxiv.org/abs/2305.07759>
3. llama2.c. URL: <https://github.com/karpathy/llama2.c>
4. LLM-ПОК Демо. URL: <https://github.com/SergeyStaroletov/LLM-ПОК>