

Triton

Kernel Triton

The aim of Triton is to provide an open-source environment to write fast code at higher productivity than CUDA, but also with higher flexibility than other existing DSLs.

<https://github.com/openai/triton>

<https://triton-lang.org/main/index.html>

Diagram illustrating the Triton compiler workflow:

- Input:** A **kernel** (represented by a 1x10 grid of boxes) and a **Triton** component (represented by a 1x10 grid of boxes).
- Process:** The **kernel** and **Triton** component are combined into a **Triton compiler** (represented by a 1x10 grid of boxes).
- Output:** The **Triton compiler** produces a **Triton** component (represented by a 1x10 grid of boxes) and a **Triton** component (represented by a 1x10 grid of boxes).

Triton [] [] [] [] [] ---Block-wise [] [] Block [] [] [] [] [] Block [] [] [] Triton compiler

[] [] [] [] [] [] [] [] Block [] [] [] [] [] [] [] [] Triton compiler [] [] []

Passes Triton NV GPU kernel
 coalescing pipeline/prefetch
 shared memory bank-conflict swizzling

Revision #1

Created 11 January 2025 09:46:28 by Colin

Updated 12 January 2025 06:35:20 by Colin