

# triton\_benchmarks

## Snapshot of your Triton benchmarks

Kernel	Work-load (1024 paths × 1 000 steps)	Avg launch time	Speed-up vs. CURAND
<b>fused_sde_stub</b> (RNG + cumsum only)	100 repeats	<b>25.97 <math>\mu</math>s</b>	<b><math>\approx 19 \times</math></b>
<b>fused_sde_control</b> (full feedback / decay logic)	100 repeats	<b>30.2 <math>\mu</math>s</b>	<b><math>\approx 16 \times</math></b>
<b>Baseline CURAND C++ kernel</b>	100 repeats	<b>0.496 ms</b>	—

*That's  $\sim 29$  ns per path-step for the control kernel on an A100 40 GB (Torch 2.7.1, Triton 3.3.1) .*

## What those numbers mean in practice

1. **GPU saturation is already solid** – a single A100 can crank  $\sim 33$  M path-steps / ms.

*One million 10 k-step paths would arrive in  $\approx 300$  ms, well below a sub-second inference budget.*

2. **Logic overhead is tiny** – adding your stochastic-control math costs  $\sim 4 \mu$ s ( $\approx 14$  %) over the bare RNG+cumsum. That confirms the extra branching and `tl.where` statements aren't a bottleneck yet.

3. **Why Triton beats CURAND** – the fused kernel keeps RNG, maths and stores in-register; the CURAND version burns cycles on separate memory passes and kernel-launch latency.

## Quick validation against BICEP + new risk dial

Pipeline stage	Latency slice	Budget after Triton win
----------------	---------------	-------------------------

Triton SDE + ENN cell	$\leq 0.04$ ms	✓ leaves ~20 ms for GNN-limit dial, Fusion $\alpha$ and router
GNN encoder (PyG on 1024-node contradiction graph)	~1.2 ms	still fits
Fusion $\alpha$ + order routing	~17 ms*	end-to-end $\approx 19$ ms

- Measured in earlier Colab runs on same GPU.

So even with the new **limit/colimit micro-risk dial** replacing the single Ising score, you remain safely under the 25 ms end-to-end target.

## Head-room & micro-optimisations

Idea	Expected gain	Effort
<b>Counter-based Philox RNG</b> (tl.philox) to eliminate per-step seed math	1-2 $\mu$ s	2 lines
<b>Shared-memory scratch for acc</b> followed by one final global store	2-3 $\mu$ s on A100	6 loc
<b>Grid-strided loop</b> (tl.program_id(axis=0) + ...) to process $\geq 2$ paths / block	Keeps occupancy high when $n_{\text{paths}} < \text{SM count}$	<10 loc
<b>Autotune block size</b> via Triton's meta['autotune']	+5-8 % throughput across GPUs	drop-in decorator

These refinements would push the control kernel toward **~22–24  $\mu$ s** for the same workload – but they're nice-to-have rather than required.

## Real-world impact for trading

- **Intraday Monte-Carlo refresh:** you could regenerate a **100 k-path, 5 k-step** Brownian cloud (~3.5 GB of output) in  $\approx 15$  ms, enabling path-level uncertainty updates every price tick without stalling the book.
- **Risk dial synergy:** the ultra-low SDE latency lets your new **limit/colimit micro-risk gauge** query fresh path statistics (vol-of-vol, excursion probabilities) in the same window it evaluates spread and depth.

- **Edge in options:** For vol surfaces you can afford to Monte-Carlo Greeks on-demand rather than rely on cached grids – handy when contradiction shocks blow out implied vol.
- 

## Bottom line

Your Triton kernels already deliver a **16–19× speed-up** over the CURAND baseline, shaving the stochastic-sampling slice down to **<0.05 ms**. That keeps the full Fusion  $\alpha$  + BICEP + limit-dial stack comfortably inside a real-time (<25 ms) trading loop, with a bit of head-room left for future bells and whistles.