

Errorless Compute For Scalable AI

James Douglas Boyd
Founder and CEO/CTO, SciSci Research

The Age of Low-Precision

AI companies are scaling models and compute because of [scaling laws](#), which suggest that model error decreases as one increases data, parameters, and compute (measured in floating-point operations per second, or FLOPS). AI Compute scales by getting more FLOPS out of GPUs. Since 2012, the greatest source of GPU FLOPS growth has been [number representation](#) for lower-precision. Indeed, NVIDIA offers GPUs that do [8-bit](#) or even [4-bit](#) floating-point arithmetic. This is necessary due to the floating-point trade-off between precision and speed.

Low-Precision Won't Continue to Scale

Looking forward, it seems unlikely that low-precision compute will continue to scale AI. The benefits are diminishing, and the risks are raised as precision is lowered. Epoch AI forecasts only a 2× [FLOP increase](#) by 2030 from 8-bit floating point adoption, a rather marginal advancement. On the

other hand, as precision is decreased further, the risks become serious, requiring either complicated workarounds or limits on lowering precision. These include [fewer effective parameters](#), reduced with [weight accuracy](#), and unstable training due to [gradient biases](#).

Beyond Floating-Point

The way to increase operations per second is not to go after FLOPS at all; floating-point is a poor format for representing numbers and doing arithmetic. SciSci has a number representation scheme of its own, based not on floating-point arithmetic, but on exact p -adic arithmetic. SciSci is building the first fast and errorless AI accelerator chip, the exact processing unit (EPU), to perform exact arithmetic without any rounding error, and faster than floating-point, using this alternative number representation. The EPU will give perfect-precision results and 15× speed gains *per operation per ALU*. SciSci aims to end the tradeoff between precision and speed for more scalable AI.