



Probabilistic computing using noisy devices

SAND2024-11633C

The authors acknowledge financial support from the DOE Office of Science (ASCR / BES) for our Microelectronics Co-Design project COINLFIPS. SNL is managed and operated by NTESS under DOE NNSA contract DE-NA0003525. This paper describes technical results and analysis. Any subjective views or opinions that might be expressed in the paper do not necessarily represent the views of the U.S. Department of Energy or the United States Government

Shashank Misra¹, **Christopher R. Allemang**¹, Christopher D. Arose¹, **Brady G. Taylor**², Andre Dubovskiy³, Ahmed Sidi El Valli³, Laura Rehm³, Andrew Haas³, Andrew D. Kent³, Leslie C. Bland⁴, Suma G. Cardwell¹, **J. Darby Smith**¹, J. Bradley Aimone¹

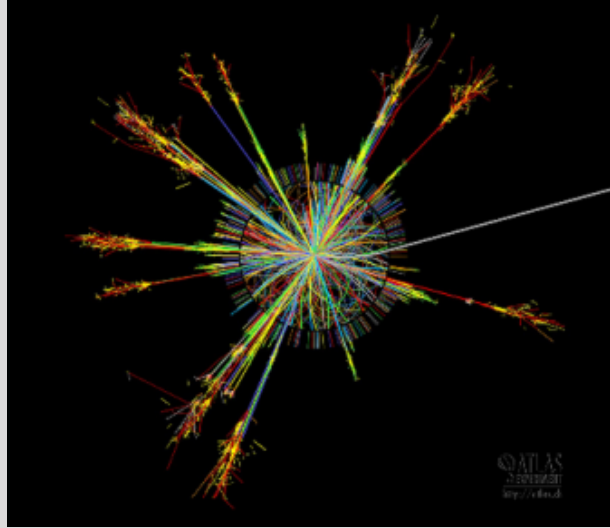
¹Sandia National Laboratories ²Duke University ³New York University ⁴Temple University

Statistical computations and the hardware lottery



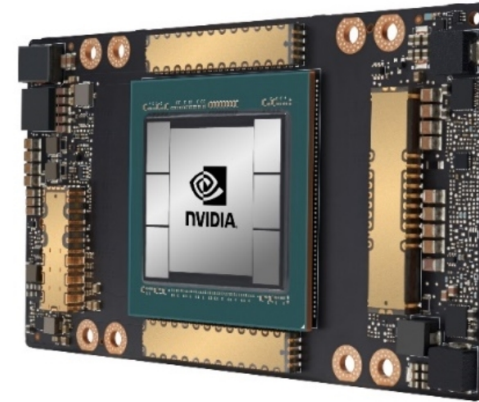
Artificial Intelligence

Bayesian neural networks are appealing yet often computationally intractable

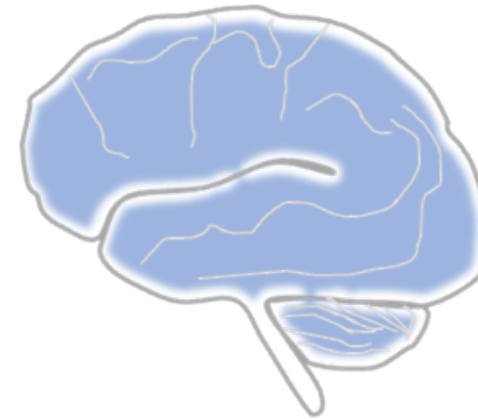


Modeling and Simulation

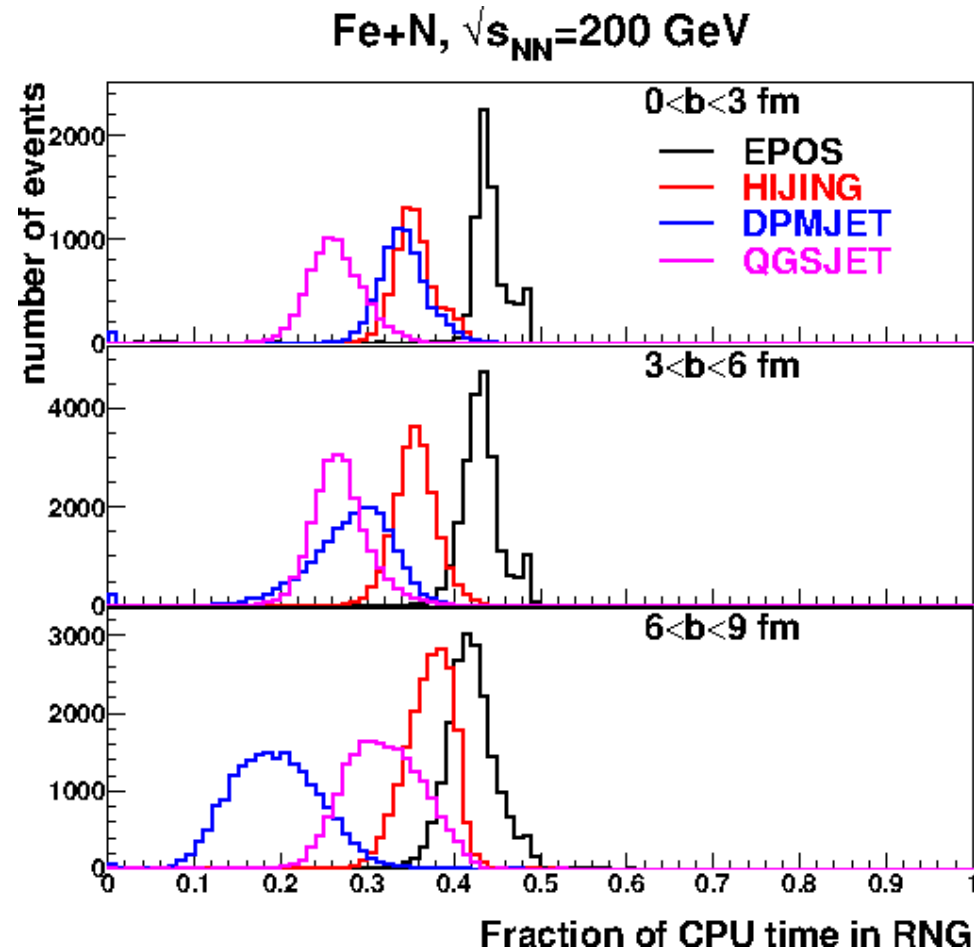
Many applications are inherently stochastic in their physics and are best modeled using probabilistic methods



~400 W
~ 10^{13} - 10^{14} FLOPS
Fully deterministic



~20 W
~ 10^{15} events / second
Fully stochastic



Misra, Adv. Mater. (2022)

Pythia - event generator used here to simulate particle showers from high energy cosmic rays

Profile amount of time used by the pseudo-random number generator (PRNG)

Uniform random sample converted to sample other distributions using elementary functions

Some calculations consume random numbers faster than they can be produced

How can we use noisy devices?

Potentially three orders of magnitude efficiency moving from pseudo random number generator (PRNG - software) to a true random number generator (TRNG - hardware)...

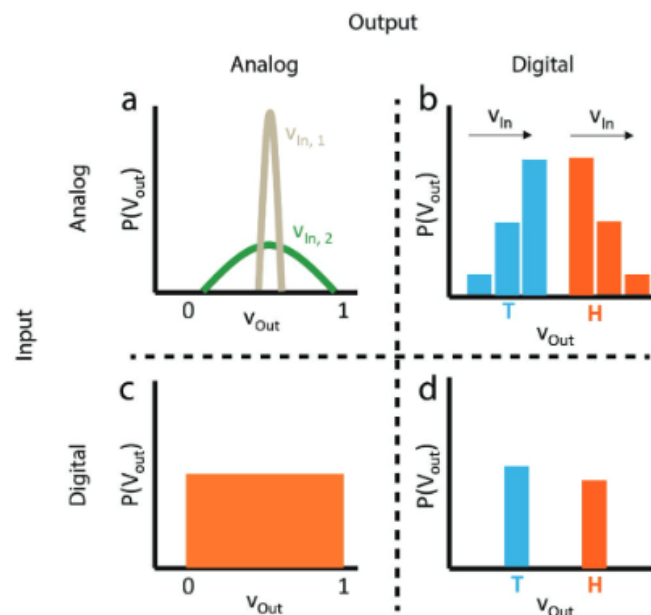
- PRNGs: ~ 1 nJ
- TRNG (MTJ, TD): < 1 pJ

Djupdal, CARRV (2023)

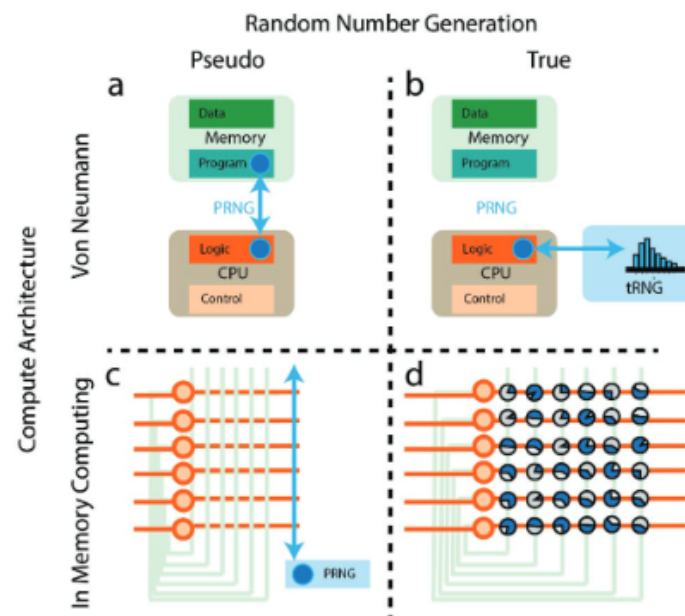
A. Shukla, IEEE ISQED (2023)

... but unclear how to use TRNGs in practice.

Fair vs. weighted?



Accelerator vs. integrated?



Misra, Adv. Mater. (2023)

Misra & Aimone, IEEE EMD (2023)

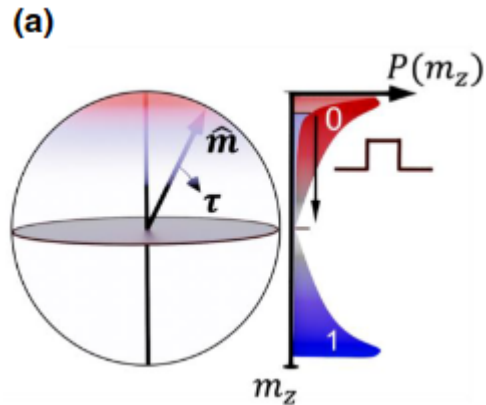
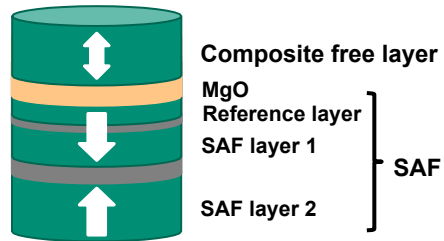
1. Hardware bitstreams
2. Sampling distributions
3. System implications



Coinflip device - a random bitstream generator

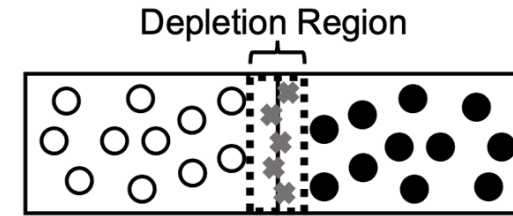


Stochastic magnetic actuated random transducer (SMART)

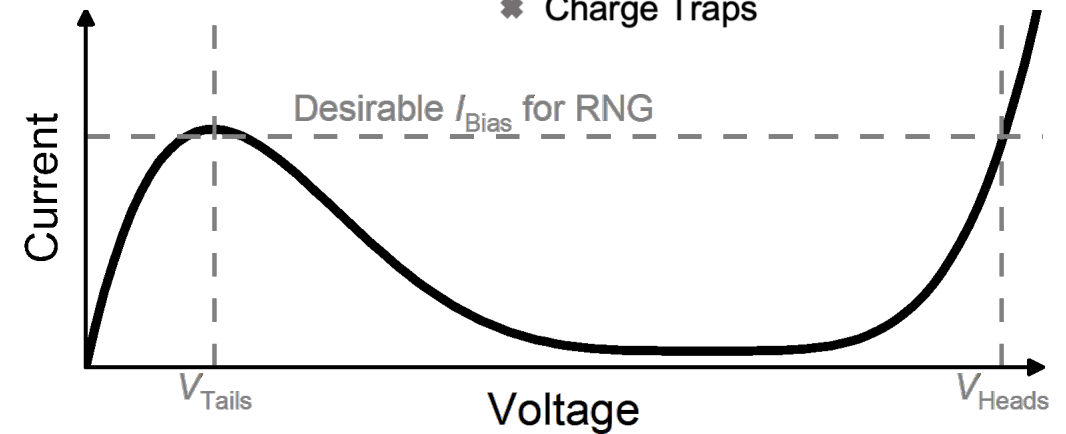


Rehm, Phys. Rev. Appl. (2023)

Tunnel diode



- Electrons
- Holes
- ✱ Charge Traps



Why?

- Output states are well-resolved
- Input pulse controls probability of 'high' output
- Easy to understand simple case with uniform sampling using fair coin

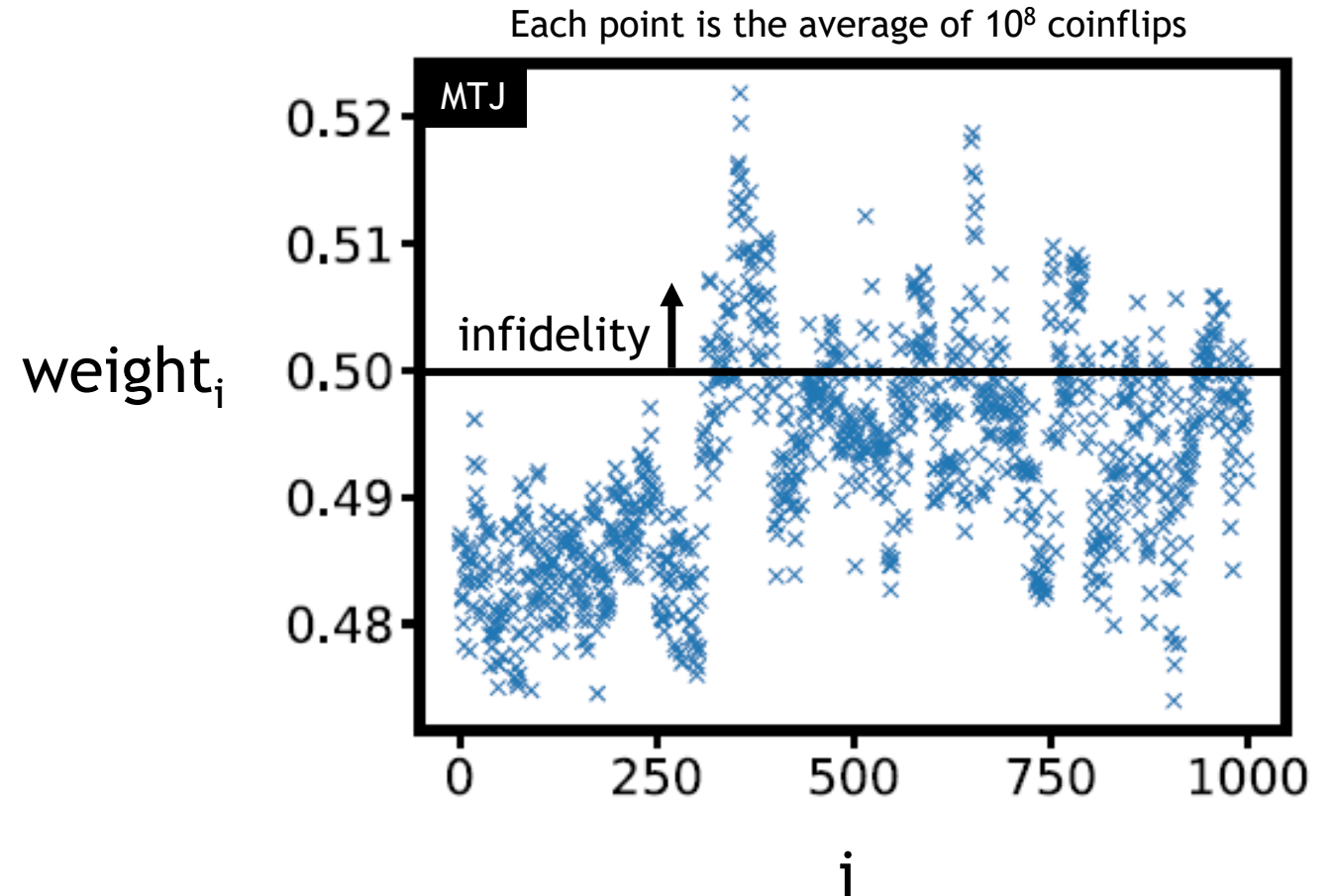
How do we evaluate the bitstream (1)?



NIST tests are useful litmus test, but provide little insight into sampling

How well can we set the weight?

Define infidelity $\delta_i = w_i - 0.5$



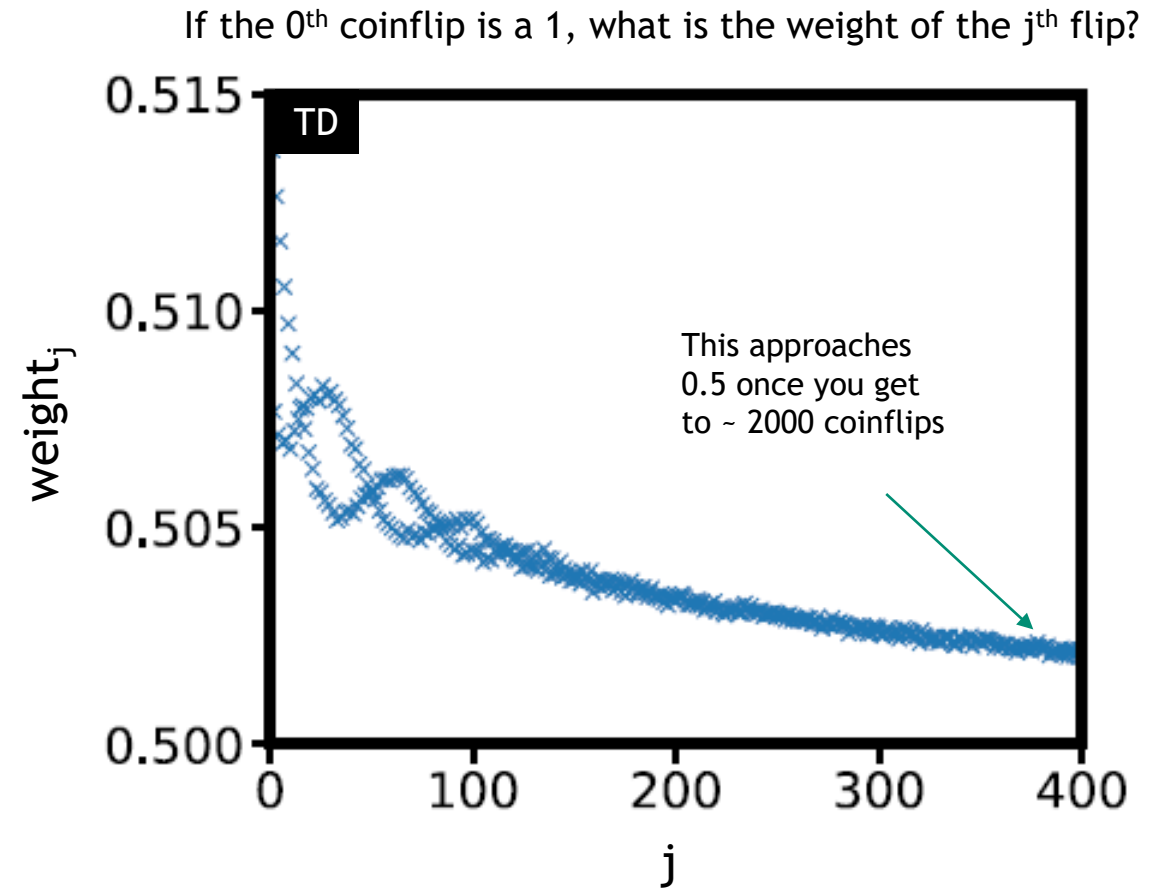
Infidelity drifts with external factors - e.g. temperature

How do we evaluate the bitstream (2)?



Are the coinflips independent of one another?

Define dependence $\varepsilon = w_1 - 0.5$



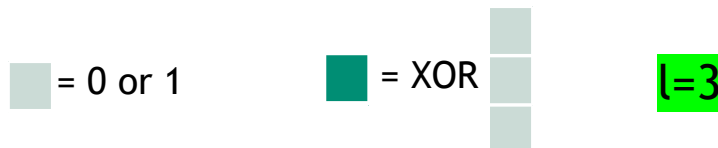
Dependence can be intrinsic (heating from last pulse), or extrinsic (pickup)

Can we improve accuracy?

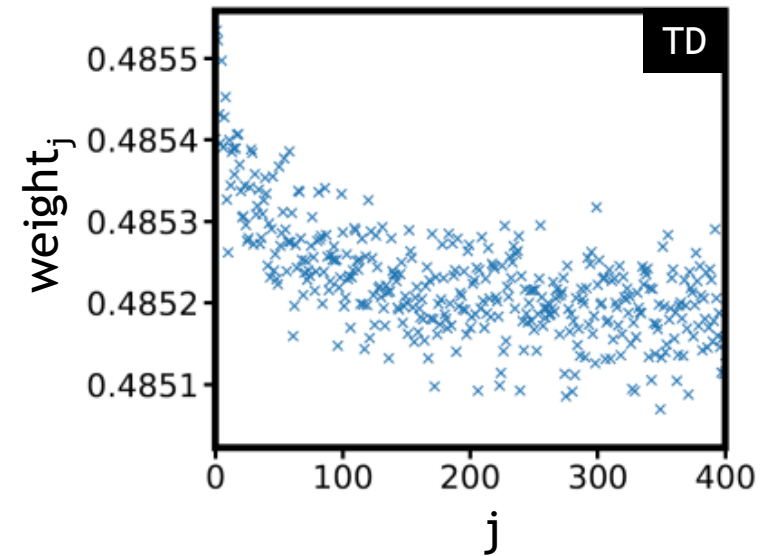
infidelity δ
dependence ε

First	Second	Raw	XOR2	XOR3
1	1	$\frac{1}{4} - \frac{\varepsilon}{2} + \frac{\delta}{2}$	$\frac{1}{4} + \frac{\varepsilon}{2} - \delta^2$	$\frac{1}{4} - 2\varepsilon\delta$
0	1	$\frac{1}{4} + \frac{\varepsilon}{2} + \frac{\delta}{2}$	$\frac{1}{4} + \frac{\varepsilon}{2} - \delta^2$	$\frac{1}{4} - 2\varepsilon\delta$
0	0	$\frac{1}{4} - \frac{\varepsilon}{2} - \frac{\delta}{2}$	$\frac{1}{4} - \frac{\varepsilon}{2} + \delta^2$	$\frac{1}{4} + 2\varepsilon\delta$
1	0	$\frac{1}{4} + \frac{\varepsilon}{2} - \frac{\delta}{2}$	$\frac{1}{4} - \frac{\varepsilon}{2} + \delta^2$	$\frac{1}{4} + 2\varepsilon\delta$

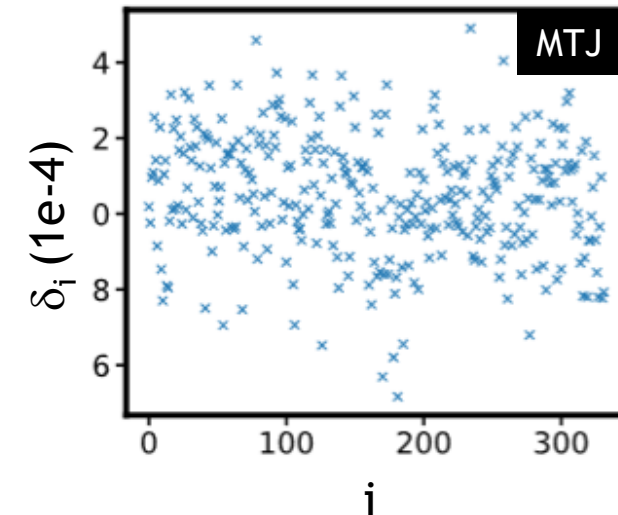
Logical exclusive or of consecutive bits:
XOR2 vs. XOR3



Dependence (XOR2)



Weight drift (XOR3)

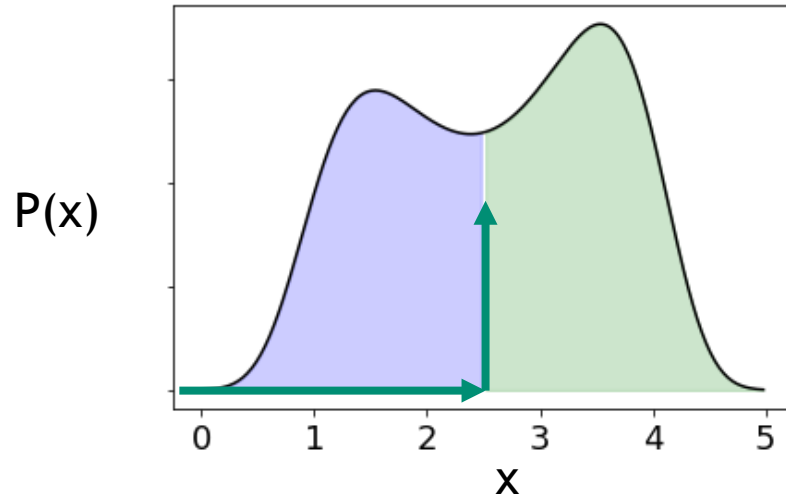


Outline

1. Hardware bitstreams
2. Sampling distributions
3. System implications



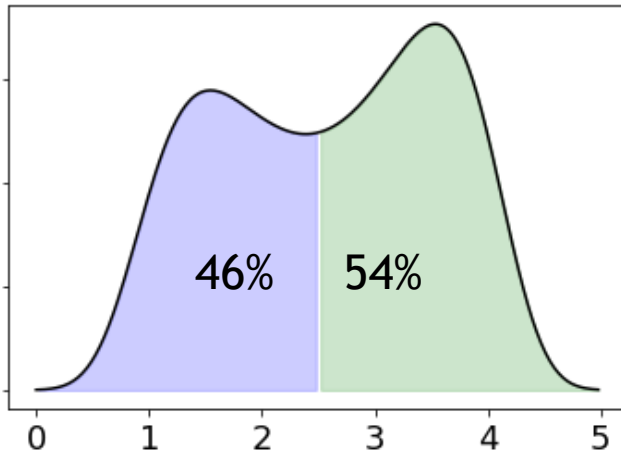
Sampling



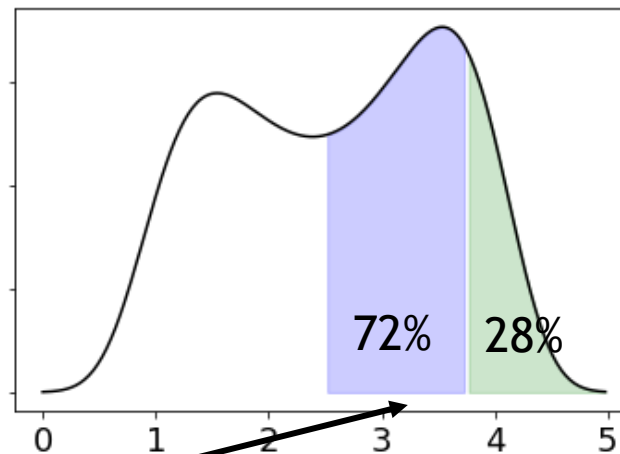
Rejection sampling:

- Random number between 0 and 5 $\rightarrow x$
- Random number between 0 and 1 \rightarrow accept/reject

Top half or bottom half?

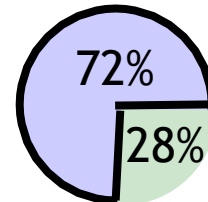
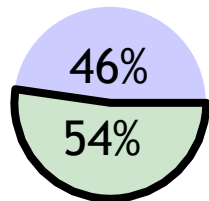


Top quarter or 3rd quarter?



Sampling tree

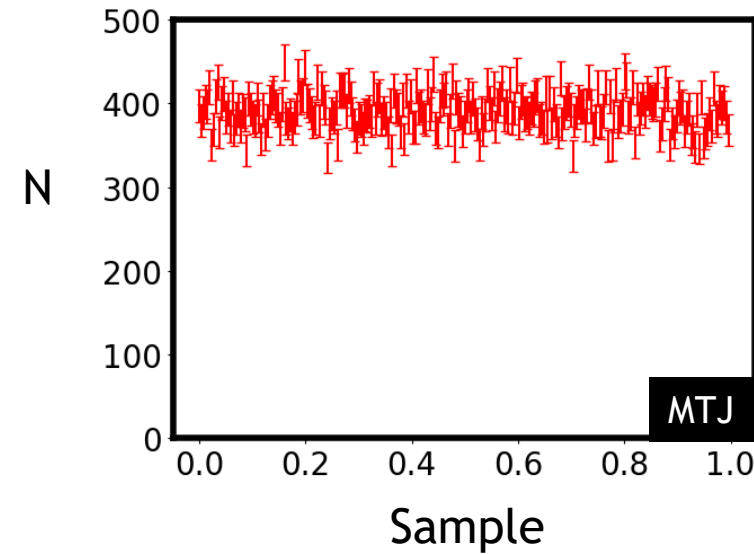
- Weighted decision determines successive bits
- Gryzka, Periodica Mathematica Hungarica (2021)*



Sampling a uniform distribution



Discretized uniform random sample



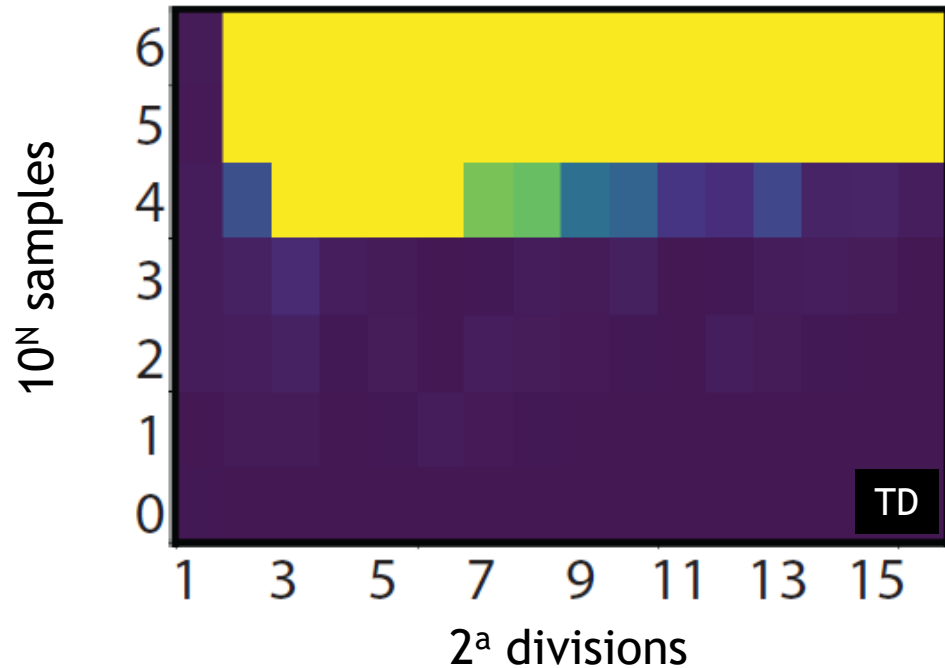
χ^2 statistics tests whether the sample and target distributions are distinguishable

Statistical uncertainty of $\sqrt{N_i}$
How you sum χ^2 is important.

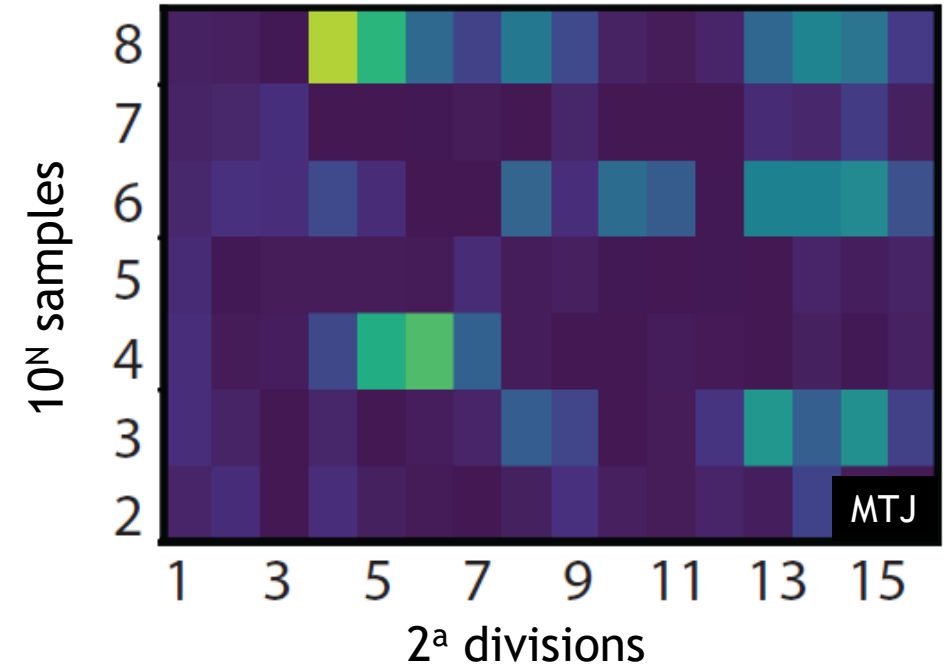
δ and ε impact sampling a uniform distribution



How much does ε matter?

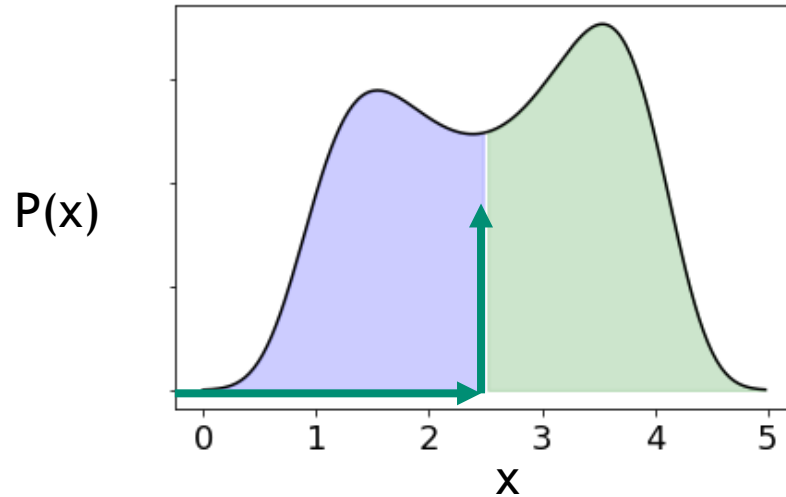


XOR3 improves sampling



Heuristic: $N \max(\delta, \varepsilon)^2 \sim 1$

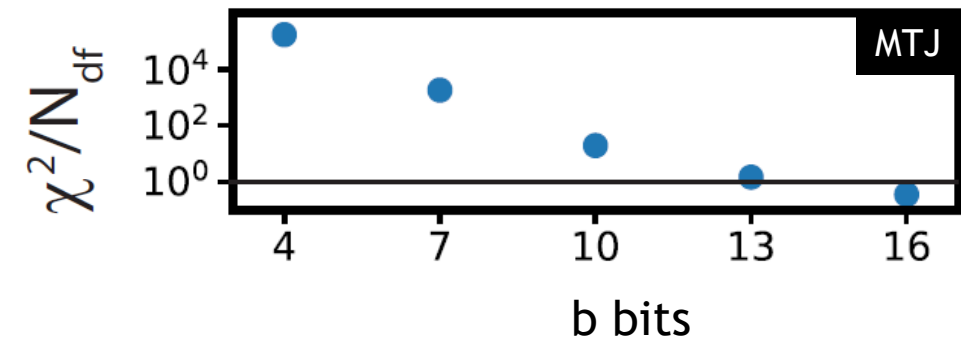
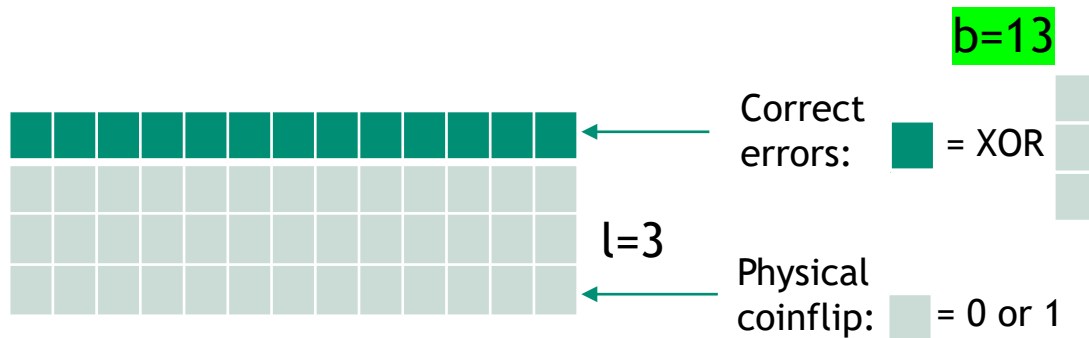
Rejection sampling



Problem: Say we want 10^8 samples - requires $\delta, \varepsilon \sim 10^{-4}$

Rejection sampling:

- Random number between 0 and 5 $\rightarrow x$ (32 bits)
- Random number between 0 and 1 \rightarrow accept/reject (b bits)

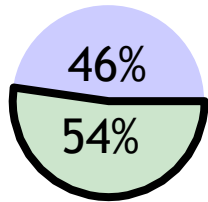
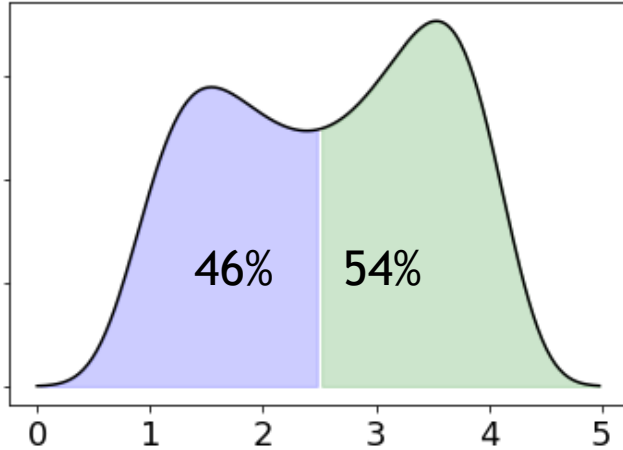


Heuristic: $N \max(1/2^b)^2 \sim 1$

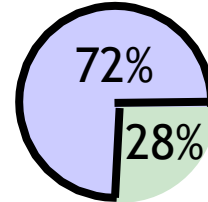
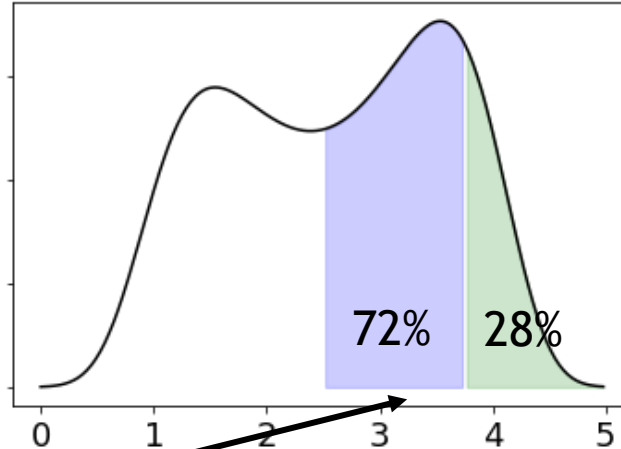
Sampling tree



Top half or bottom half?



Top quarter or 3rd quarter?

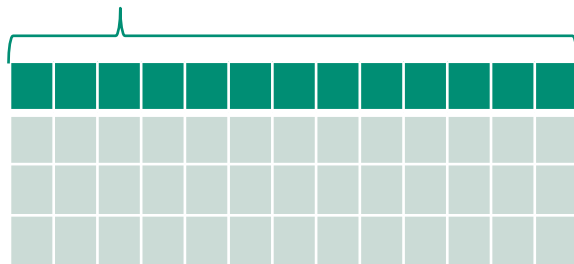




...

Problem: Say we want 10^8 samples - requires $\delta, \epsilon \sim 10^{-4}$

Impractical for a weighted coinflip device.

Logical weighted coinflip

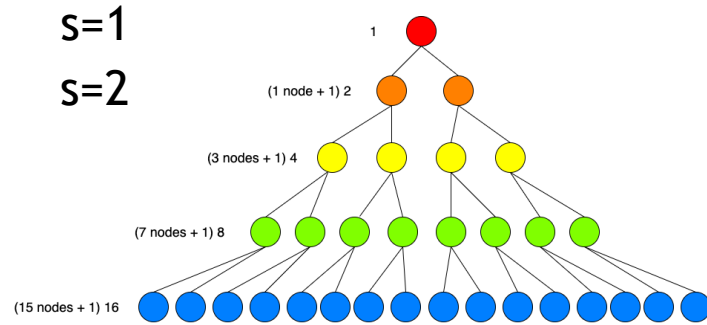


Correct errors:  = XOR 

Physical coinflip:  = 0 or 1

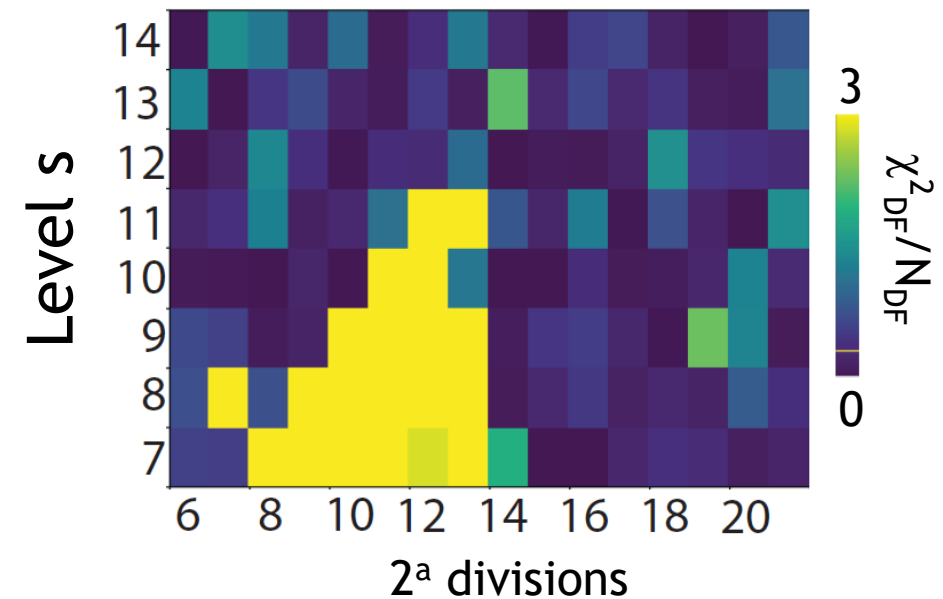
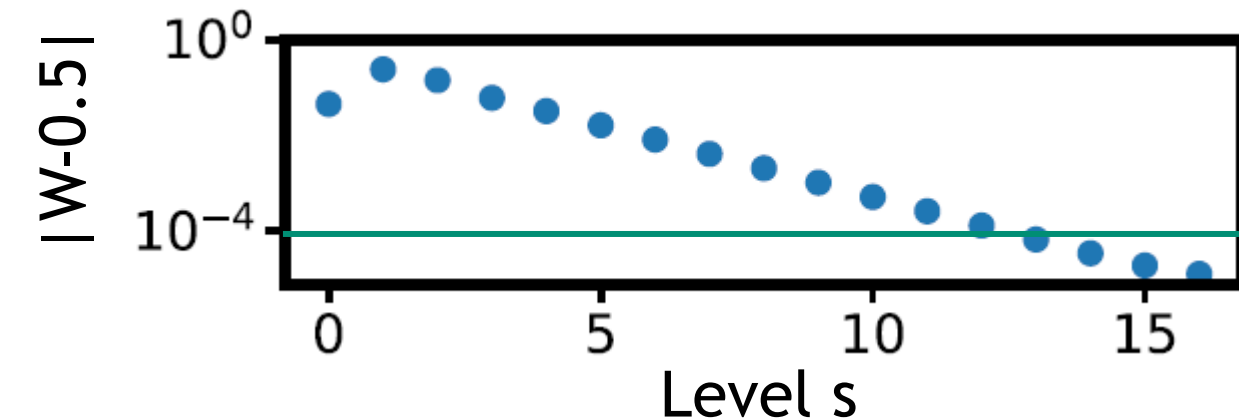
Solution: use fair coins to draw a uniform random sample with 13 bits of precision

Cutoff sampling tree for efficiency

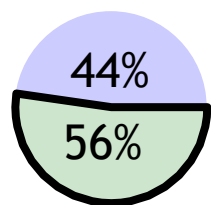
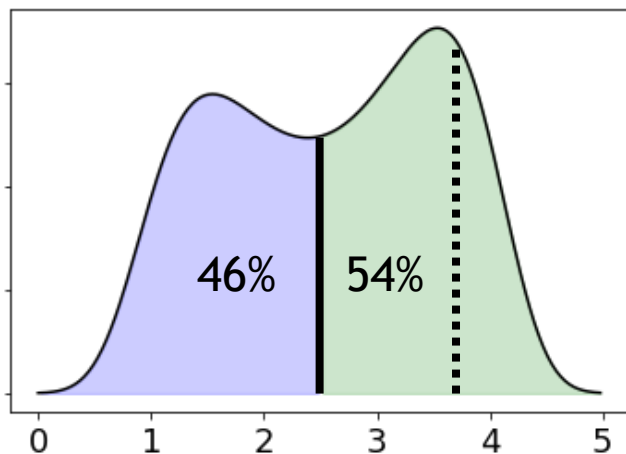


Need to store $2^{32}-1$ sixteen bit weights in memory, but nearly all of the weights are 0.5 (after $s=12$)

Only need sampling tree for top 12 bits - remaining bits can be uniform random sample



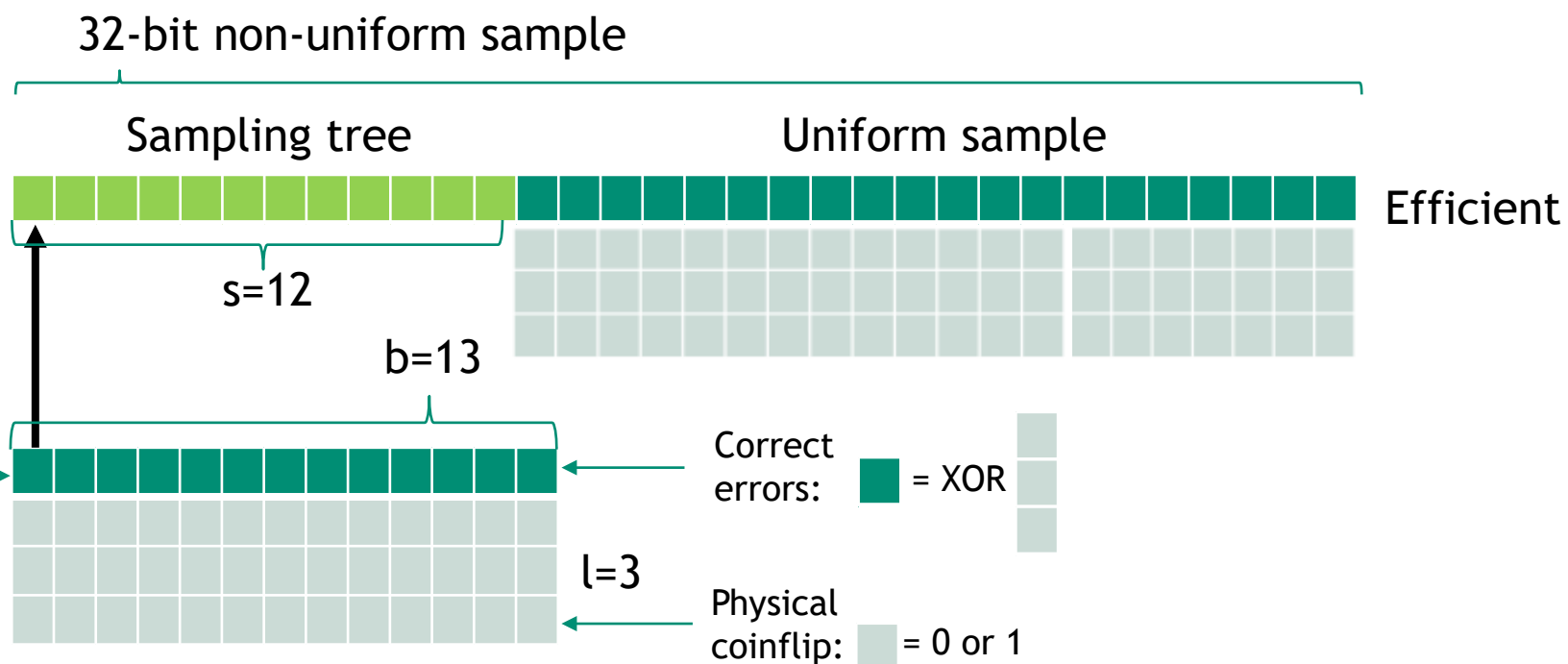
Scheme



Vs.

Uniform random number

Accurate



Outline

1. Hardware bitstreams
2. Sampling distributions
3. System implications



What device is better?



Component	Description	Technology	Power consumption	# Needed	Reference
V_{in}	Square wave generator using operational transconductance amplifiers	130 nm CMOS	457 μ W	Potentially 1	[4]
Current Mirror	High-compliance ultra-high output resistance with positive shunt feedback	130 nm CMOS	570 μ W @ 96 μ A	1/TD	[5]
Tunnel Diode	Three-terminal silicon surface junction tunneling device	1 μ m CMOS	~1 nW in “tails” ~1.75 nW in “heads”	1/TD	[6]
Comparator	Low-noise, low-power, dynamic latched comparator using cascode	130 nm CMOS	100 μ W	1/TD	[7]
XOR	2-input pass-transistor XOR	130 nm CMOS	231 pW	Potentially 1/RN	[8]

At a few ns timing, device (aJ - fJ) < circuit (pJ) < INT (nJ)

- Device is irrelevant
- Have flexibility at the circuit level to incorporate, e.g., analog feedback
- It is not that integer operations (INT) need to get more efficient... you just want to minimize them

[4] I. A. A. Al-Darkazly *et al.*, IEEE Access, vol. 4, pp. 3169–3181, 2016, doi: 10.1109/ACCESS.2016.2557843.

[5] M. H. Maghami *et al.*, International Journal of Circuit Theory and Applications, vol. 43, no. 12, pp. 1935–1952, 2015, doi: 10.1002/cta.2049.

[6] J. Koga *et al.*, Appl. Phys. Lett., vol. 70, no. 16, pp. 2138–2140, Apr. 1997, doi: 10.1063/1.118970.

[7] A. M. Maghraby *et al.*, 2020 12th International Conference on Electrical Engineering (ICEENG), Jul. 2020, pp. 335–338. doi: 10.1109/ICEENG45378.2020.9171746.

[8] N. Ahmad *et al.*, Active and Passive Electronic Components, vol. 2013, p. e148518, Mar. 2013, doi: 10.1155/2013/148518.

Is TRNG better than PRNG? What sampling scheme is better?

Minimize integer operations (INT)

Uniform distribution

Non-uniform distribution

PRNG

TRNG

Rejection

Tree

10 INT

96 coinflips
2 XOR

2 RNG
100 INT acceptance
1 conditional
2x executed on average

13 RNG
12 conditionals
12 cache access

	TRNG	PRNG
Rejection	202	242
Tree	24	154

Misra, in preparation

Codesign *matters*

A vision for probabilistic computation

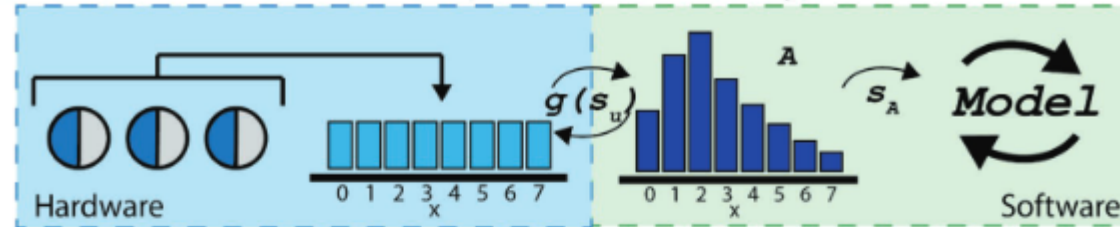


Sampling is expensive (time, energy)
Sampling < 50% of CPU time

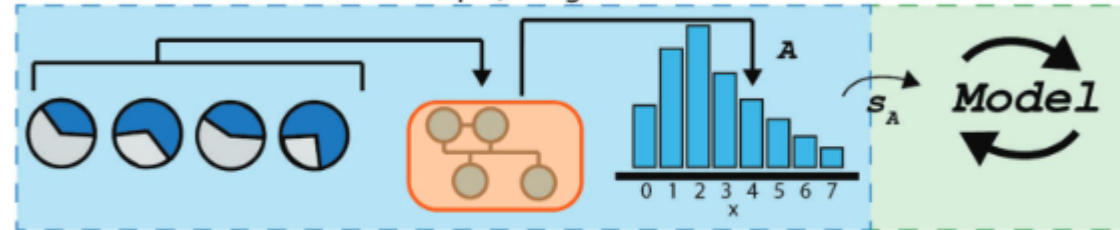
Goal: Make sampling cheap
(devices, circuits)

Move more of computational
burden from deterministic
computations into sampling
(new algorithms,
architecture, etc.)

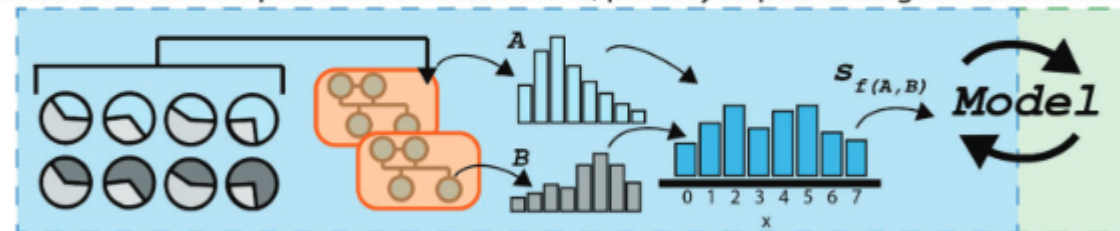
Level 0: Hardware draw uniform sample, convert to desired type in software



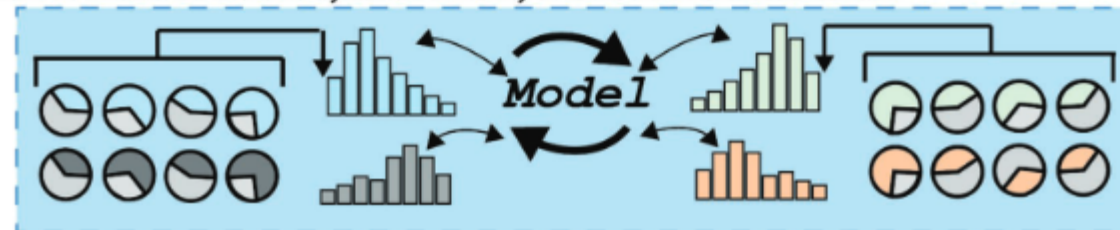
Level 1: Hardware draw desired sample, integrate into software model



Level 2: Perform computation on distributions, partially implementing model in hardware



Level 3: Model stochastic systems entirely within stochastic hardware



Today

Future