Arjun Earthperson PhD Candidate, PRA Group

NC STATE





- Motivation

Evolving Hardware Landscape



Industry Response to Emergent AI/ML Workloads

Heavy Investment in Data-Parallel Hardware

- GPUs, tensor cores provide high throughput for integer operations.
 - Apple M4 Neural Engine: ≈ 38 TOPS (int8).
 - Nvidia RTX 4090: ≈ 1000 TOPS.
- Current-gen consumer hardware already supports specialized ops (Intel AMX).

Designed for Inference on Massive Models

- $\blacksquare \approx 10^9$ parameters on mobile devices (e.g., Gemma 2B).
- $\blacksquare \approx 10^{12}$ parameters on HPC/cloud (LLaMa 4 at 400 Billion).

Comparatively,

- **Largest (public) PRA models:** $\approx 10^6$ parameters (Generic PWR).
- However, PRA model quantification is deeply recursive.



(日)

Birds' Eye View

Research Contribution



Overview and Highlights I

High-Throughput Boolean Logic Evaluation (aka Eval Query)

- Simultaneous evaluation of *all* intermediate gates, success, and failure paths.
- Relax coherence constraints: arbitrary graph shapes, with NOT, or any other combination of gates permitted.
- Boolean logic manipulation not strictly needed.
- Vectorized bitwise hardware ops for logical primitives (AND, OR, XOR, etc.)
- Specialized treatment of k/n logic, without expansion.
- Concurrent use of all available compute GPUs, multicore CPUs.



Birds' Eye View

Research Contribution



Overview and Highlights II

Probability Estimation using Eval Query (aka Expectation Query)

- Estimate probabilities for *all* events in the PRA model simultaneously.
- No need to compute the minimal cut sets or prime implicants.
- Streamable: Solve the entire PRA model iteratively, approx 0.3 seconds per iteration, regardless of model size. More iterations needed for complex models.



ヘロト 人間 トイヨト ヘヨト 三日



A Working Example: One Initiating Event, Three Fault Trees, Six Basic Events, Five End States



Variable	Expression
X	$(A B') \bullet (A' (B \bullet C'))$
Y	$C \bullet (D E)'$
Ζ	$kn[(A \bullet C), (D \bullet E), F']$

Table: Unsimplified Boolean expression foreach Top Event

Small, but non-trivial structure:

- Basic events are shared.
- Some gate outputs are negated.

イロト 不留 とうほとう ほど

• Event Z is a (k=2) of n=3 gate.

A Working Example: One Initiating Event, Three Fault Trees, Six Basic Events, Five End States

Compile a Directed Acylic Graph (DAG) from Logic Model



Start with an Arbitrary Topological Ordering:

- Aim for succinctness.
- prefer HW-native ops.
- e.g. don't expand k/n.

・ロト ・雪 ト ・ ヨ ト





PRA Models as Probabilistic Circuits

└─ Knowledge Compilation and Queries

Querying the Compiled Knowledge Graph

The Simplest Type of Query: Eval(G)

- Set the inputs [on/off].
- Observe the outputs [on/off].
- Can be used as a building block for an embedding ML model.

8/24

But just how fast is it?

NC STATE







PRA Models as Probabilistic Circuits

Knowledge Compilation and Queries

Eval Query on the Compiled Knowledge Graph

Eval Query Performance on GPUs:

- Latency: 200-300 *ms* per graph pass.
- Throughput: VRAM bound (see plot).
- Benchmarked on Nvidia GTX 1660 [6GB].
- Graph sizes: from ≈ 50 to ≈ 2000 nodes.
- Evals: from 16M to 1B per node per pass.

Q: Are these enough samples to estimate the Expectation Query?







@ANS

Knowledge Compilation and Queries

Estimator for the Expected Value (i.e., Probability)

- A Boolean function *F*(**x**) can be viewed as an indicator function: *F*(**x**) ∈ {0,1}.
 The event {*F*(**X**) = 1} has probability E[*F*(**X**)].
- Monte Carlo estimator:

$$\widehat{P}_N = \frac{1}{N} \sum_{i=1}^N F(\mathbf{x}^{(i)}),$$

where each $\mathbf{x}^{(i)}$ is a random draw from the input distribution.

By the Law of Large Numbers,

$$\lim_{N\to\infty} \widehat{P}_N = \Pr[F(\mathbf{X}) = 1], \text{ almost surely.}$$

Error decreases at rate $O(1/\sqrt{N})$, analyzed via the Central Limit Theorem.

NC STATE



@ΔNS

L Knowledge Compilation and Queries

Monte Carlo Sampling

- Rather than summing or bounding all combinations of failures, *simulate* random draws of **X**.
- Each Monte Carlo iteration:
 - **1** Sample $x_1, x_2, \ldots, x_n \stackrel{\text{i.i.d.}}{\sim} \prod p(x_i)$.
 - **2** Evaluate the Boolean function $F(\mathbf{x})$ (cost is just logical gate evaluation).
 - **3** Collect whether $F(\mathbf{x}) = 1$ (failure) or 0 (success).
- Repeating for many samples {**x**⁽¹⁾,...,**x**^(N)} yields a *sample average* estimate of the probability.
- Benefits:
 - Bypasses explicit inclusion-exclusion expansions.
 - Straightforward to parallelize (evaluate each draw in separate threads or blocks).



@ANS

・ロト ・ 母 ト ・ ヨ ト ・ ヨ ト

Back to Working Example: One Initiating Event, Three Fault Trees, Six Basic Events, Five End States









Preliminary Case Study: Aralia Dataset



@ANS

Preliminary Case Study

- Preliminary Case Study: Aralia Dataset



Overview: Aralia Dataset

- **Dataset Composition:** The Aralia collection consists of 43 trees.
- Diverse Problem Sizes: Small trees (e.g. 25–32 basic events) through mid-sized models with over 1,500 BEs.
- Wide Probability Range: Top-event probabilities spanning from rare events near 10⁻¹³ to fairly likely failures with probability above 0.7.
- Model Variability: Some trees are primarily AND/OR, others incorporate more advanced gates (K/N, XOR, NOT), providing thorough coverage of typical (and atypical) fault tree logic structures.



イロト 不得 トイラト イラト 二日

- Preliminary Case Study: Aralia Dataset
 - Benchmarking Procedure



Benchmarking Setup: Hardware and Environment

Target Hardware:

GPU: NVIDIA[®] GeForce GTX 1660 SUPER (6 GB GDDR6, 1,408 CUDA cores).
 CPU: Intel[®] CoreTM i7-10700 (2.90 GHz, turbo-boost, hyperthreading).

Software Stack:

- SYCL-based (AdaptiveCpp/HipSYCL), with LLVM-IR JIT for kernel compilation.
- Compiler optimization at -03 for efficient code generation.
- Repeated runs (5+) to mitigate transient variations.
- Measured Time: Includes entire wall-clock duration, from host-device transfers and JIT compilation to final result collection.



Preliminary Case Study: Aralia Dataset

Benchmarking Procedure



@ANS

Monte Carlo Execution and Implementation

- **Objective:** Compute TOP event probabilities for all 43 trees.
- Sampling Strategy:
 - Single pass per fault tree, generating as many samples as fit in 6 GB GPU memory.
 - 128-bit Philox4x32x10 pseudo-random number generator, parallel threads.

Bit-Packing Optimization:

- Each group of 64 Monte Carlo outcomes stored in a single 64-bit word.
- Enables vectorized instructions (e.g. popcount) and reduces memory I/O.

Data Types:

- Tallies in 64-bit integers.
- Probability accumulations in double precision (64-bit float).

æ

©ANS

Preliminary Case Study: Aralia Dataset

NC STATE

- Accuracy Benchmark: Relative error (Log-probability), Data-Parallel Monte Carlo vs Min-Cut Upper Bound and Rare-Event Approximation



18/24

Preliminary Case Study: Aralia Dataset

NC STATE

Performance Benchmark (Memory Consumption): Sampled Bits Per Event Per Iteration







©ANS

— Outlook

Limitations:

- Brute-force/naive Monte Carlo struggles when sampling rare-events.
 - Implement importance sampling: WIP.
- Brute-force/naive Monte Carlo is a poor strategy for sampling correlated events.

Next Steps:

Benchmark on larger (G-PWR, G-MHTGR) models.

Future Work:

- Embedding models from Knowledge Graph, for semantic representation.
- Gradient computation on Knowledge Graph.
- Minimal cut set computation from Knowledge Graph.



・ロト ・ 母 ト ・ ヨ ト ・ ヨ ト



Outlook



The End

NC STATE



(日)

-Outlook

Knowledge Compilation



Hierarchy of compiled target languages. Blue nodes represent canonical forms.



Acronym	Full form
NNF	Negation Normal Form
XAG	XOR-And-Inverter Graph
AIG	And-Inverter Graph
ANF/RNF	Algebraic/Ring Normal Form
f-NNF	Flat Negation Normal Form
DNNF	Decomposable Negation Normal Form
d-NNF	Deterministic Negation Normal Form
FPRM	Fixed Polarity Reed-Muller
CNF	Conjunctive Normal Form
DNF	Disjunctive Normal Form
s-DNNF	Smooth/Structured Decomposable Negation Normal
d-DNNF	Deterministic Decomposable Negation Normal Form
sd-DNNF	Smooth/Structured Deterministic Decomposable Neg
PPRM	Positive Polarity Reed-Muller
PI	Prime Implicate
IP	Prime Implicant
BCF	Blake Canonical Form
EPI	Essential Prime Implicate
EIP	Essential Prime Implicant
BDD	Binary Decision Diagram
f-BDD	Free/Read-Once Binary Decision Diagram
OBDD	Ordered Binary Decision Diagram 🚊 , 🚊 🥠 🤈 📀
SDD	Sentential Decision Diagram
RoBDD	Reduced Ordered Binary Decision Diagram

21/24

Preliminary Case Study: Aralia Dataset

Convergence Trends

NC STATE



@ANS

Convergence over 1000 iterations, Aralia das9204



Input Data





э

Preliminary Case Study: Aralia Dataset

L Input Data

23	edfpa14b	311	290	70	-	-	-	105,955,422	2.95620E-01
24	edfpa14o	311	173	42	-	-	-	105,927,244	2.97057E-01
25	edfpa14p	124	101	42	-	-	-	415,500	8.07059E-02
26	edfpa14q	311	194	55	-	-	-	105,950,670	2.95905E-01
27	edfpa14r	106	132	55	-	-	-	380,412	2.09977E-02
28	edfpa15b	283	249	61	-	-	-	2,910,473	3.62737E-01
29	edfpa15o	283	138	33	-	-	-	2,906,753	3.62956E-01
30	edfpa15p	276	324	33	-	-	-	27,870	7.36302E-02
31	edfpa15q	283	158	45	-	-	-	2,910,473	3.62737E-01
32	edfpa15r	88	110	45	-	-	-	26,549	1.89750E-02
33	elf9601	145	242	97	-	-	-	151,348	9.66291E-02
34	ftr10	175	94	26	-	-	-	305	4.48677E-01
35	isp9601	143	104	25	1	-	-	276,785	5.71245E-02
36	isp9602	116	122	26	-	-	-	5,197,647	1.72447E-02
37	isp9603	91	95	37	-	-	-	3,434	3.23326E-03
38	isp9604	215	132	38	-	-	-	746,574	1.42751E-01
39	isp9605	32	40	8	6	-	-	5,630	1.37171E-05
40	isp9606	89	41	14	-	-	-	1,776	5.43174E-02
41	isp9607	74	65	23	-	-	-	150,436	9.49510E-07
42	jbd9601	533	315	71	-	-	-	150,436	7.55091E-01
43	nus9601	1,567	1,622	392	47	-	-	unknown	



