# Pitfalls and learnings in performance modeling

Phil Thierry*, Cedric Andreolli, Sai Chenna , Fabrice Dupros, Sunny Gogar,
Sylvain Jubertie, Nalini Kumar, Amine Mrabet and Mariam Umar
Intel Corp.

intel.

# Introduction / Conclusion

One Simulator cannot satisfy all requirements and solve every questions

Rule number one:
- Define the objective :
    - One-shot question / project  or should we develop of full modeling Workflow
- Define the input parameters available and how to get them (are the tools ready)
- Define the output results and needed accuracy
    - If a Lifetime is not enough, revisit the hardware and software granularities
    - In any case, use uncertainties
    - And better have at least 2 or 3 ways to validate the results

Analytical models
- Any Scale. Most complexities (comms. , compute) are implicitly accounted for.
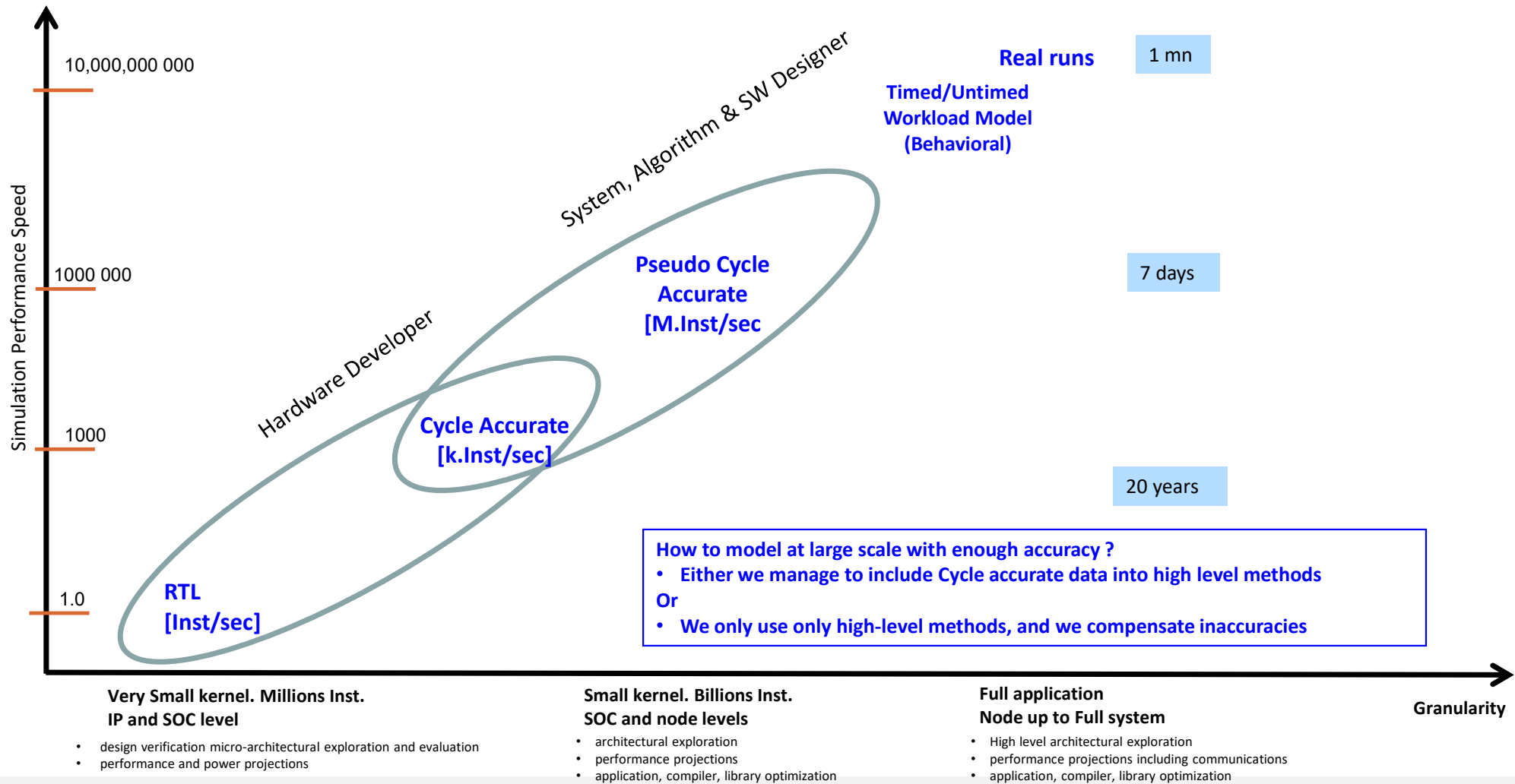
Graph-based modeling (AI-NN)
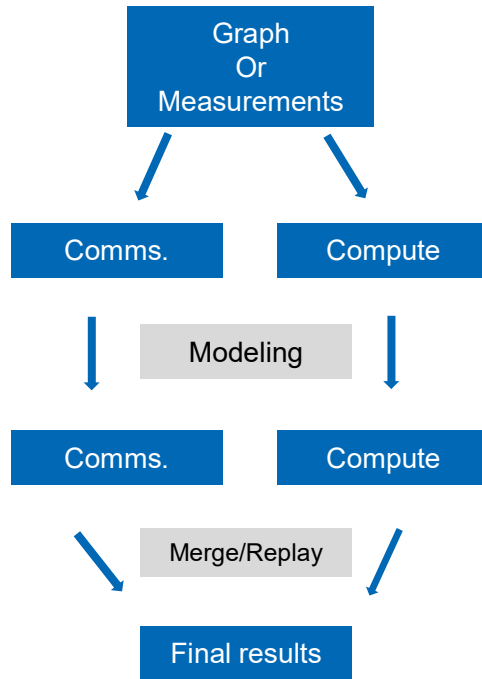- Any scale but limited to NN

Measurement-based modeling (HPC and AI)
- Single CPU/GPU core to single CPU/GPU core
- Single homogeneous/heterogeneous node to Single homogeneous /heterogeneous node
- O(x) ranks to  O(x) ranks
- O(x) ranks to  O(10x) ranks

# Accuracy and speed are orthogonal features !



**Simulation Performance Speed** (y-axis)

- 10,000,000 000
- 1000 000
- 1000
- 1.0

**Real runs**

**Timed/Untimed Workload Model (Behavioral)**

1 mn

7 days

20 years

**Pseudo Cycle Accurate [M.Inst/sec**

**Cycle Accurate [k.Inst/sec]**

**RTL [Inst/sec]**

System, Algorithm & SW Designer

Hardware Developer

**How to model at large scale with enough accuracy ?**
- **Either we manage to include Cycle accurate data into high level methods**
Or
- **We only use only high-level methods, and we compensate inaccuracies**

**Granularity** (x-axis)

**Very Small kernel. Millions Inst.**
**IP and SOC level**
- design verification micro-architectural exploration and evaluation
- performance and power projections

**Small kernel. Billions Inst.**
**SOC and node levels**
- architectural exploration
- performance projections
- application, compiler, library optimization

**Full application**
**Node up to Full system**
- High level architectural exploration
- performance projections including communications
- application, compiler, library optimization

# The right methods for the given applications

Graph
Or
Measurements

Comms.          Compute

Modeling

Comms.          Compute

Merge/Replay

Final results

Analytical modeling

HPL/ HPLMxp

HPCG-UQ

3DFFT

FD/QCD

LLM

- **Graph-based methods for AI: (mostly python)**
    "Collective communication + GEMM and convolution"
- **Measurements-based methods for HPC and AI:**
    MPI, OMP, C/C++/ Sycl/ Fortran, etc...
    Distributed and Shared memory + Any numerical schemes

    A way to mix accuracy and speed

- Parametric equation based on the numerical schemes

- Simple model based on F/B

    Fast, Any scale. Can be adapted to specific requests
    Complexity is implicitly accounted for.

# It's all about the communication & compute relation (IO too)

Application can be described as  Total_time = f [ **Comms** (+IO), **Compute** (+ serial_init) ]

Every modeling are doing  comms $\oplus$ compute*

Where $\oplus$ can be either
- An addition :  not accounting for overlap, dependencies, asynchronism
- A more complex "replay" (scheduling of comms and compute)

|  | Communications | Compute |
|---|---|---|
| profiling | APS, SMPI, SCOREP, ONNX | SDE,  Advisor, MSR (Vtune, Emon, Papi), Unitrace |
| modeling | LogP , SIMGRID, SST | Analytical , NN,  Sniper, Cy.Accurate |

Graph
Or
Measurements

Comms. | Compute

Modeling

Comms. | Compute

Merge/Replay

Final results

Then How to chose:
1. Speed :  we may quantify them as processed Inst / sec
2. Accuracy: we can define ranges (<5%  , [5-10], [10-20], [20,50],[50,100])
3. Scalability: how easy to move to another application

**Back to first rule:  Clearly define the objectives**

# AI-NN specific case

Identify volume of communication at every stages (mostly collectives)

Identify GEMM, convolution per layers (mostly size and shapes)

Extrapolate to a given number of end points using
- Analytical collective models
- Topology aware tools (SST, cycle acc., …)

Extrapolate them
Roofline extrapolation,
(pseudo) Cycle Accurate simulation

Merge communications and compute timing
- Summation with or without overlap
- Reuse original scheduling (if DAG exist)

```
                    Graph
                   /      \
            Comms.          Compute
                   Modeling
            Comms.          Compute
                 \          /
                  Merge/Replay
                  Final results
```

# Graph-based methods

3 [{"LayerNum": 3072, "LayerName": "Level3072_0", "Dependency": ["Level3072_0", "Level3106_1", "Level3156_10", "Level3230_100", "Level3340_101", "Level3450_102", "Level3560_103", "Level3634_104", "Level3708_105", "Level3818_106", "Level3928_107", "Level4038_108", "Level4148_109", "Level4258_11", "Level4332_110", "Level4442_111", "Level4516_112", "Level4590_113", "Level4700_114", "Level4810_115", "Level4920_116", "Level5030_117", "Level5140_118", "Level5250_119", "Level5324_12", "Level5398_120", "Level5448_121", "Level5522_122", "Level5596_123", "Level5670_124", "Level5744_125", "Level5818_126", "Level5892_127", "Level5942_128", "Level5992_129", "Level6066_13", "Level6140_130", "Level6214_131", "Level6288_132", "Level6362_133", "Level6436_134", "Level6510_135", "Level6560_136", "Level6634_137", "Level6744_138", "Level6854_139", "Level6964_14", "Level7038_140", "Level7148_141", "Level7258_142", "Level7368_143", "Level7442_144", "Level7516_145", "Level7626_146", "Level7736_147", "Level7846_148", "Level7956_149", "Level8066_15", "Level8116_150", "Level8226_151", "Level8300_152", "Level8374_153", "Level8484_154", "Level8594_155", "Level8704_156", "Level8814_157", "Level8924_158", "Level9034_159", "Level9108_16", "Level9158_160", "Level9232_161", "Level9342_162", "Level9452_163", "Level9562_164", "Level9672_165", "Level9782_166", "Level9892_167", "Level9966_168", "Level10040_169", "Level10150_17", "Level10224_170", "Level10334_171", "Level10444_172", "Level10554_173", "Level10664_174", "Level10774_175", "Level10848_176", "Level10922_177", "Level11032_178", "Level11142_179", "Level11252_18", "Level11326_180", "Level11436_181", "Level11546_182", "Level11656_183", "Level11730_184", "Level11780_185", "Level11854_186", "Level11928_187", "Level12002_188", "Level12076_189", "Level12150_19", "Level12224_190", "Level12298_191", "Level12348_192", "Level12382_193", "Level12432_194", "Level12482_195", "Level12532_196", "Level12582_197", "Level12632_198", "Level12682_199", "Level12716_2", "Level12766_20", "Level12840_200", "Level12890_201", "Level12964_202", "Level13038_203", "Level13112_204", "Level13186_205", "Level13260_206", "Level13334_207", "Level13384_208", "Level13434_209", "Level13508_21", "Level13582_210", "Level13656_211", "Level13730_212", "Level13804_213", "Level13878_214", "Level13952_215", "Level14002_216", "Level14052_217", "Level14126_218", "Level14200_219", "Level14274_22", "Level14348_220", "Level14422_221", "Level14496_222", "Level14570_223", "Level14620_224", "Level14670_225", "Level14744_226", "Level14818_227", "Level14892_228", "Level14966_229", "Level15040_23", "Level15090_230", "Level15164_231", "Level15214_232", "Level15264_233", "Level15338_234", "Level15412_235", "Level15486_236", "Level15560_237", "Level15634_238", "Level15708_239", "Level15758_24", "Level15808_240", "Level15858_241", "Level15932_242", "Level16006_243", "Level16080_244", "Level16154_245", "Level16228_246", "Level16302_247", "Level16352_248", "Level16386_249", "Level16436_25", "Level16510_250", "Level16560_251", "Level16610_252", "Level16660_253", "Level16710_254", "Level16760_255", "Level16794_26", "Level16868_27", "Level16942_28", "Level17016_29", "Level17090_3", "Level17140_30", "Level17214_31", "Level17264_32", "Level17314_33", "Level17388_34", "Level17462_35", "Level17536_36", "Level17610_37", "Level17684_38", "Level17758_39", "Level17808_4", "Level17858_40", "Level17908_41", "Level17982_42", "Level18056_43", "Level18130_44", "Level18204_45", "Level18278_46", "Level18352_47", "Level18402_48", "Level18452_49", "Level18526_5", "Level18576_50", "Level18650_51", "Level18724_52", "Level18798_53", "Level18872_54", "Level18946_55", "Level18996_56", "Level19030_57", "Level19080_58", "Level19130_59", "Level19180_6", "Level19230_60", "Level19280_61", "Level19330_62", "Level19380_63", "Level19414_64", "Level19464_65", "Level19538_66", "Level19612_67", "Level19686_68", "Level19760_69", "Level19834_7", "Level19868_70", "Level19942_71", "Level19992_72", "Level20066_73", "Level20176_74", "Level20286_75", "Level20396_76", "Level20506_77", "Level20616_78", "Level20726_79", "Level20800_8", "Level20850_80", "Level20924_81", "Level21034_82", "Level21144_83", "Level21254_84", "Level21364_85", "Level21474_86", "Level21584_87", "Level21658_88", "Level21732_89", "Level21842_9", "Level21916_90", "Level22026_91", "Level22136_92", "Level22246_93", "Level22356_94", "Level22466_95", "Level22540_96", "Level22614_97", "Level22724_98", "Level22834_99", "Level3071"], "MsgSizeInBytes": 8, "CommCycles": "0.837816", "Time": "56076.311683", "MPICall": "MPI_ALLREDUCE", "HotPhaseId": "1014", "Rank": 0, "comm": "-2080374781", "OPTYPE": "comm"},
4 {"LayerNum": 3073, "LayerName": "Level3073", "Dependency": ["Level3072_0"], "OPTYPE": "compute", "ComputePhaseTag": "p1014c0001i0000", "ComputeTimeMSec": "87.91339999999764", "Rank": 0},
5 {"LayerNum": 3074, "LayerName": "Level3074_0", "MsgSizeInBytes": 872616, "CommCycles": "0.028413", "Time": "56165.062899", "MPICall": "MPI_IRECV", "HotPhaseId": "1014", "Rank": 0, "comm": "-1006632959", "src": "1", "tag": "1", "OPTYPE": "comm", "Dependency": ["Level3073"]},

**Format conversion**     **+ Hdw config and scheduling**

Dependency flow : "Scheduling"

12   "nodes": [
13     {
14       "label": "Level3072_0",
15       "typedef": "Layer",
16       "data": {
17         "Layer": {
18           "Layer Name": "Level3072_0",
19           "Layer Index": 3072,
20           "l_pass": "fwd",
21           "fwd_pass_msg_size": 8.0,
22           "fwd_pass_comp_cycles": 0.0,
23           "comms_time_fwd_cycles": 0.0,
24           "comms_scaleout_time_fwd_cycles": 0.0,
25           "source": -1,
26           "destination": -1,
27           "GPU_tiles": 2,
28           "GPU_cards": 8,
29           "CPU_nodes": 16,
30           "fwd_pass_collective": "MPI_Allreduce"
31         }
32       },
33       "id": 0

- Msg size
- Computes cycles
- Comms Time scaleup (T2T+Scaleup)
- Comms Time scaleout
- Source Rank/Dest Rank

377566   "links": [
377567     {
377568       "source": "Level3072_0",
377569       "target": "Level3072_0",
377570       "key": "key1"
377571     },
377572     {
377573       "source": "Level3106_1",
377574       "target": "Level3072_0",
377575       "key": "key2"
377576     },
377577     {
377578       "source": "Level3156_10",
377579       "target": "Level3072_0",
377580       "key": "key3"
377581     },

# Traces format Standardization

- Using Python model APIs : Custom dev for the full workload

  - What we do so far. Some advantages but many disadvantages

- Using ONNX model format (OpenNeural Network exchange)

- Using Chakra Execution trace

ONNX exchanges models between various frameworks,
Chakra's exchanges execution traces between different teams.

Choosing the right "API" is extremely important to scale the methods to various application and model

*see Sridharan et al. 2023. "Chakra.: Advancing Performance Benchmarking and Co-design using Standardized Execution Traces."*

# HPC and AI cases

Profile volume of communication
Who is talking to who, imbalance
Aggregated (APS, vampyr) or full tracing
(SIMGRID, SCOREP)

Extrapolate to a given number of end points
using
- BW/lat scaling for aggregated comms
- Topology aware tools (SST, SIMGRID, cycle acc., ...)

Merge communications and compute timing
- Summation with or without overlap
- Reuse original scheduling (if DAG exist)

Profile the compute (sampling, hdw counters)
- Aggregated
- Per function
- Tracing

Extrapolate them with
Roofline extrapolation,
(pseudo) Cycle Accurate simulation

```
                        Graph

        Comms.                    Compute

                     Modeling

        Comms.                    Compute

                   Merge/Replay

                   Final results
```

# Amdahl's law inflection point

Modeling when and why the scalability is failing down remains a huge challenge, for both strong and weak scaling.  Not to mention IO.

We need it for Interconnection/Topologies whatif analysis and pathfindings.

With analytical modeling and graph-based AI methods, we can mimic the communication and compute overlap.
Doable with HPL, 3dFD, 3dFFT, LQCD, AI LLM (i.e deterministic applications "easy" to formalize).

For more complex applications, it remains impossible to profile on O(x) ranks and extrapolate to O(10x) ranks (i.e extrapolation of scalability).

The only way would be a full tracing (time independent) of those O(10x) ranks to capture the real communications patterns and compute.

Too expensive with SIMGRID. Same with SCORE-P. Traces quickly becomes far too large

Our current tentative is with SCORE-P using a trace compressor (phases detection) that works with regular application.

This remains  "WIP". Not to mention the Ultimate: Modeling productions runs with several apps running at same time.

# SST features for network simulation

Ember – State machine model to generate network traffic:

- **Features:**
  - Provides a collection of HPC comm patterns and MPI Collectives ("motifs")
- **COMET* extensions:**
  - OneCCL Allreduce Algos (includes topo-aware pipelined Allreduce)
  - Enable overlap compute and communication events for AI Modeling

Firefly – Interface b/w network driver (Ember) and router (Merlin):

- **Features:**
  - Provides detailed NIC model and packetization and byte movement engine
  - Implementation of MPI communication protocols
- **Limitations:**
  - No Non-blocking Collectives

Merlin – Low-level, high-radix router component:

- **Features:**
  - Models flit-level movement, physical routing and delivery of packets
  - Provides readily available topologies with static and dynamic routing (congestion-aware)
- **COMET extensions:**
  - Included additional topologies: Hierarchical networks, Polarfly/Polarstar,
- **Limitations:**
  - No collective/switch offload
  - IP specific Congestion control routing schemes

*COMET – Intel extension of SST.



"Application-less" patterns modeling for various (realistic) topologies compared to analytical (theoretical) model.

This is just a scale-up example with 24 endpoints in a heterogeneous node.

Analytical collective models are often too optimistic, and this effect keep growing with scale-out.

# SST various usages



SST as a SCORE-P backend for Communications and compute replay for HPC and AI "measurements-based" methods

- Extract relevant information from traces
  - System Counters (perf)
  - Communication Phases: P2P blocks , Collective Operations
- Combine trace events from different independent ranks
- Merging P2P operations into MPI Phases
- Detecting Compute Phases (gaps)
- Detecting Cyclic Phases / Loops -> Reduce Tracefile size and modeling time/complexity
- Feeding application communication phases into SST



Application-less patterns modeling for various topologies



As an alternative to Comms analytical modeling for AI-NN apps.

*COMET – Intel extension of SST.

# SIMGRID modeling. Topo, latency and BW on real applications



SIMGRID modeling. Normalized elapsed time for varying bandwidth/latency parameters.
Left-panel (flat topology) – Right-panel (2-level topology). Colors represent various BW [0.1 : 4 TB/s] and the series denotes various latencies

**Advantages** :  Full real application and test cases. No instrumentation. SIMGRID is an active opensource project.

As for SST, topologies can be added. Can be coupled with compute models. Handle comms & comp dependencies.

**Limitations**:  Can scale up to O(3) nodes but getting slow after that.

# Compute phases extrapolation from measurements

A CPU+GPU performance extrapolations with a (extended) Roofline Model



- Based on measurement to collect flops and memory traffic (Emon for CPU , Unitrace on GPU)

- Traces are realigned and metrics are merged per time sample (usually 10ms)

- For each sample, we determine if performance is bounded by CPU or GPU

- Efficiency is extracted on current platforms and projected on new platforms

# Analytical modeling

# Analytical HPCG-UQ



Analytical HPCG model

HPCG = 0.201 (F/B HPCG) → **Hand counted**
         * Peak_bw* eff_vs peak → **100%READ BW**
         * efficiency HPCG → **What HPCG can extract**
         * efficiency comms → **How does HPCG scale**

As those few parameters are uncertain, just add distributions

```
perf_dist = dev_num*
arith_int_dist
•   peak_bw * eff_100R_dist
•   *eff_HPCG_dist
•   * eff_scala_dist
```

**Initial distribution depends on *a priori* knowledge**

**Even with a perfect model, using single hdw/sftw parameters gives a single data point, which cannot be accurate, or be very lucky !**

**HPCG-UQ modeling on 4608 PVC [mean 1264.58 TF/s], vs HPCG Measurements done on 768 Aurora nodes [1259.94 TF/s] .**
**In Green are the HPCG and Stream PDF.In Blue are the uncertainties on input parameters. (Phil)**

# HPL Analytical modeling @ 1.012 EF/s



HPL can be approximated with XLS very quicky but the Analytical model can simulate compute and communications panel after panel.
( with random noise to mimic machine noises,  CPU and GPU efficiency, dimensions, BWs)

Limitations : No topology. Need to mimic the real implementation. Hard to model all parameters behavior (P, Q, NB variations)

SIMGRID can do it too but with limited number of nodes.

# Analytic LLM



Structure of transformer-based LLMs (megatron- GPT, Llama...)

Transformer based LLMs: are built using transformer decoder blocks. With each block comprising two sub-components:

- Multi-head Attention
- MLP (Multi-layer perception)

We calculate the number of floating-point operation for each sub-components:

- Hidden size, The number of attention heads, the sequence size
- The training batch size, micro batch size
- Number of transformer blocks



| | Analytical model | PVC RUN | Error |
|---|---|---|---|
| Num flop (Tflop) | 3433 | 3037,334798 | 13% |
| Elapsed time (second) | 153 | 160,8311873 | -5% |
| Throughput (tokens/second) | 2342,48366 | 2228,423516 | 5% |

Lamma training 1,46B : Comparison of the analytical model and the results of the runs on 1 tile PVC

**Advantages** :  The result is accurate, which can be attributed to the compute-bound nature of this case test

**Limitations** : Handling various parallelisms (projection at scale) and memory impact (size, bw, lat).  WIP in the community

# Conclusion

# Conclusions

- There is no single tool to model all hardware and applications details at every scale

- The definition of the objectives is the most important starting point

  - Clear Definition of input / output, i.e what questions we want to solve

    - a new u-arch / a SOC / a system platform , switch, IO, memory tiers, power, TCO

  - Definition of the "application" scale, probably the easiest

    - a few instructions / a kernel /a proxy application / a full application

- Accuracy and speed remains orthogonal

  - There are orders of magnitude between each level.

  - Somehow difficult to merge granularity but  propagating uncertainties is mandatory.

  - Full tracing remains very difficult.

- Collection tools are critical.  Investing in opensource community and standardization is crucial.

- Modeling Scientific computing  and AI applications is different, but AI remains a HPC application

intel.

# Notices and Disclaimers