**h:**

nts

$10^0$

$10^{-1}$

$10^{-2}$

$10^{-3}$

$10^{-4}$

$10^{-5}$

$10^{-6}$

Runtime / Event [s]

N Events [#]

$10^2$  $10^3$  $10^4$  $10^5$  $10^6$  $10^7$  $10^8$

Graph: On-board CPU (Cold)
Graph: On-board GPU (Cold)
Graph: Server GPU (Cold)
Graph: Server GPU (Warm)
Typical Executable

**COMPLEX**

Optimiser & Placeholder

$p_x^1$  $p_y^1$

$p_x^2$  $p_y^2$

**Vectorised Neutrino**

**Computin**

mentum:
ics

nass

atics:
olution

$W$

$l$

$\nu$

Martin Erdmann, Peter Fackeldey

IDT-UM Collab

Legend:

Array    Analytical Operator
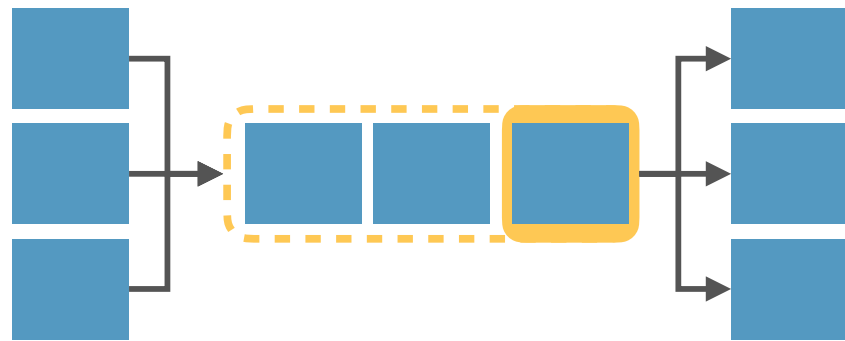
Tensor    Fitting Operator

Results

- Fast turn-around times of HEP analyses are driver of scientific insight

- Traditional analyses already O(weeks)

- Data increased in HL-LHC by x20

- Future analyse must be: Faster & More resource efficient

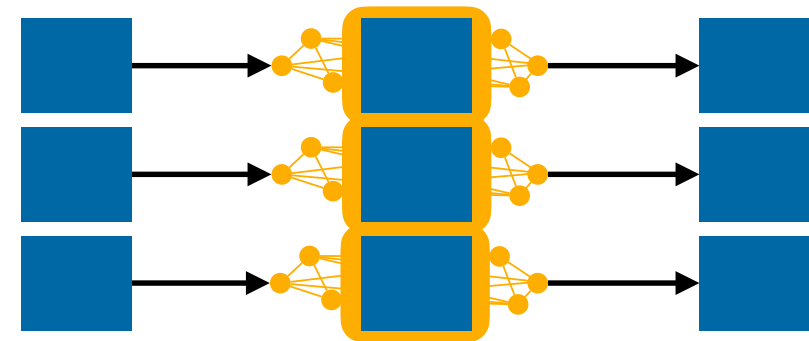➔ **Requires re-thinking of analysis computing!**

# Computations and Graphs

- Fast and efficient analyses can be realised with vectorised computations
  moving to vectorised computations can be realised with vectorised computations
  challenging to vectorise:
- Already used by many analyses

**Typical Executable:**

Events Event

# Computing Graph

Key functionalities are part of the graph:
- Analytical operations
- ...ser & placeholders for fitting
  ...perations:

Can be s...

E.g. N...

## Event loop



**Graph:**

ny events

Load one event

Evaluate expressions

Store results

Repeat

- Problem: Some computations challenging to vectorise!
- E.g. Neutrino reconstruction (next page)

no momentum:
nematics

## Vectorised

...th 100 iterations
...ptive optimisation (Adam)
...evaluated during execution



REAL

$p_Z^1$

$p_Z^2$

Choice

Event Kinematics

COMPLEX

Optimiser & Placeholder

$p_x^1$

$p_x^2$

$p_y^1$

$p_y^2$

Choice

$p_Z$

**Experimen**

- Reconstruction of longitudinal neutrino momentum in e.g. ttH events
- Solved by assuming: $\not{E}_T = p_{\nu,T}$, W mass
- Inputs: Lepton, $\not{E}_T$
- Two branches:
  - Real branch (h≦1): Purely analytical
  - Complex branch (h>1): Involves fitting

$$p_{\nu,z}^{1,2} = \frac{k}{p_{l,T}^2} \left( p_{l,z} \pm E_l \sqrt{1 - \underbrace{\left( \frac{p_{l,T}\, p_{\nu,T}}{k} \right)^2}_{\equiv h}} \right)$$
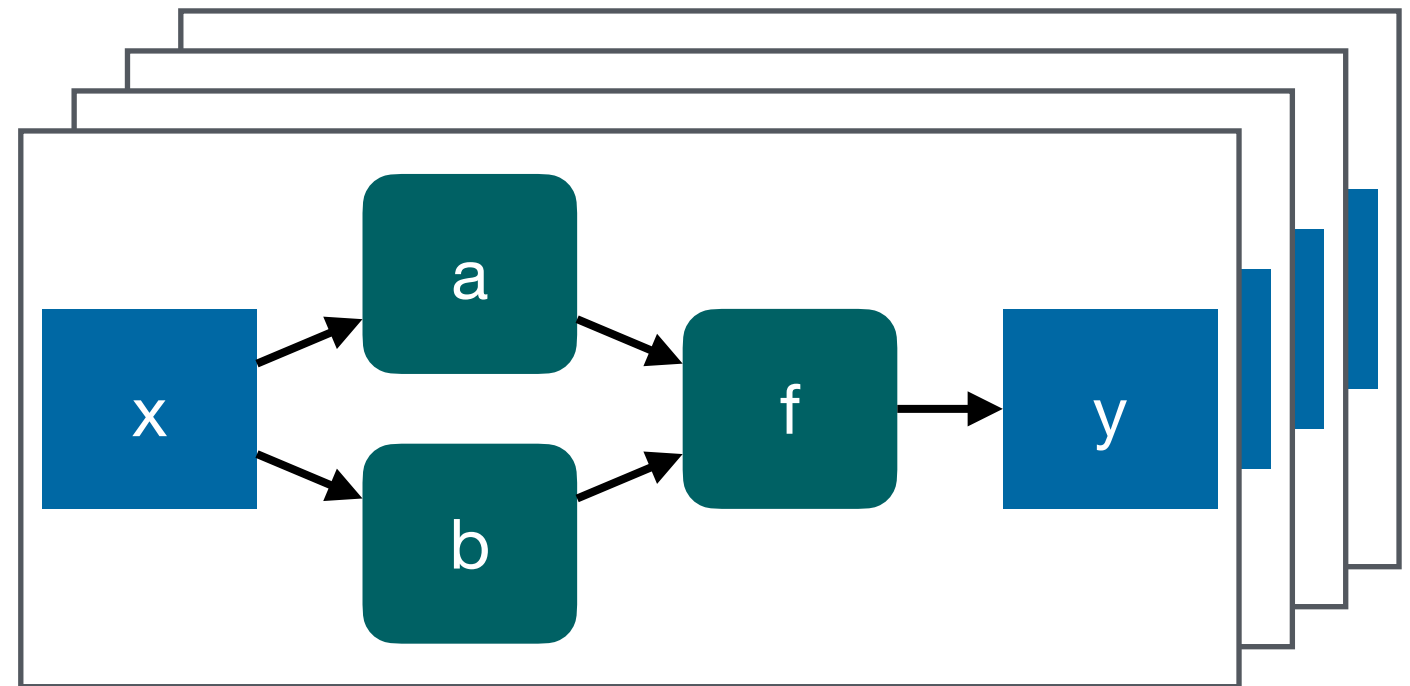
Optimal benchmark

- Multiple Inputs (Lepton, MET)
- Different behaviour on event-basis
- Stateful computation: Fitting
- Physics result can easily be verified

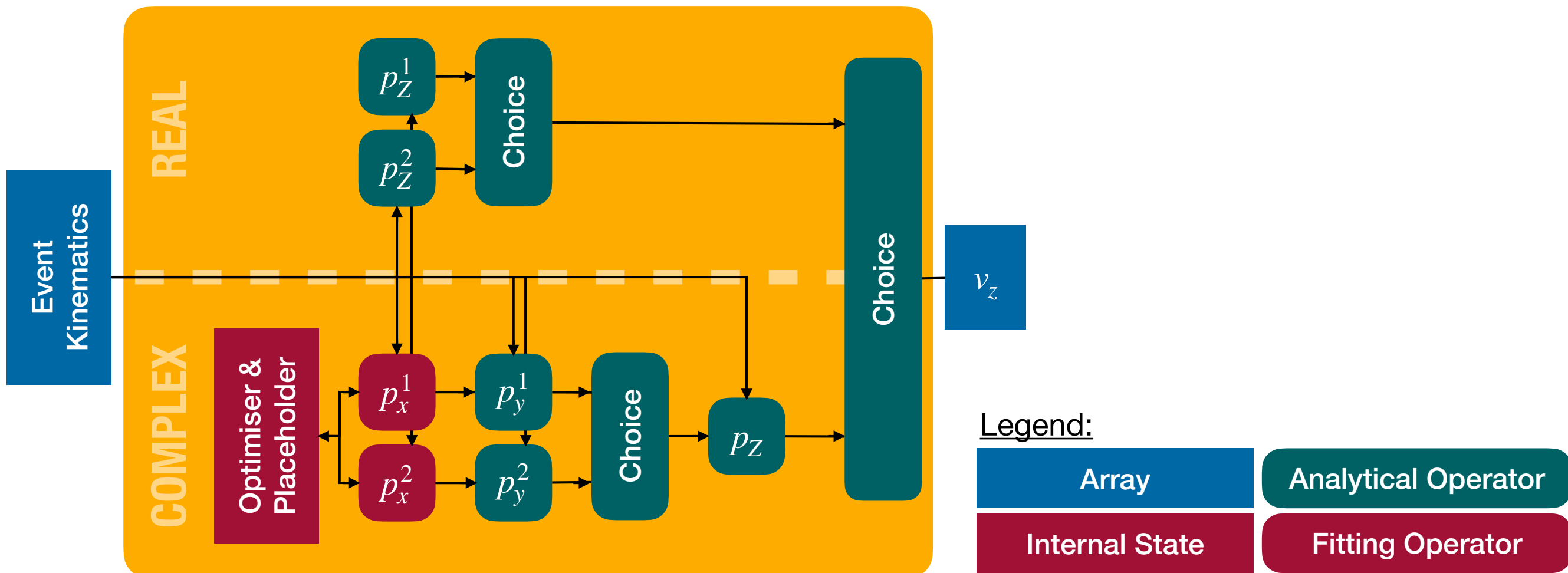➜ **Vectorise using graph computing model!**

**Computing Graph**

- Contents:
    - Nodes = Computations
    - Edges = Data flow
- Properties:
    - Directed = ⟶
    - Acyclic = no loops
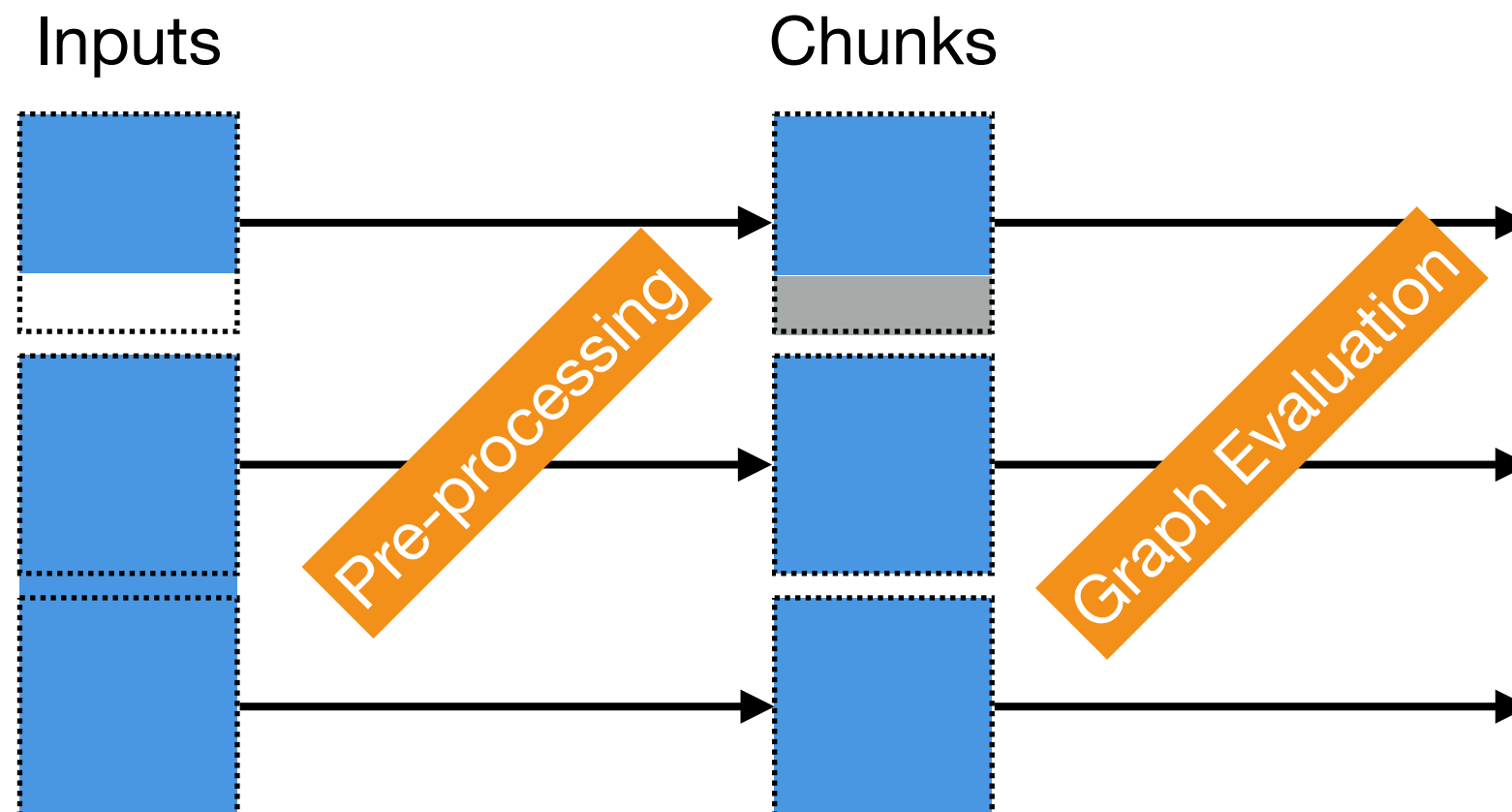
Example $y = f(a(x), b(x))$:



- Two levels of parallelism:
    - Inter processing unit:
        - Parallel units in directed acyclic graph (e.g. a & b) can run in parallel
    - Intra processing unit (SIMD):
        - If graph is same for multiple inputs ($N_{events}$)
        - Parallel execution over many events

- Two branches (real and complex)
- Fit of neutrino momentum:
  - Unrolled for-loop with 100 iterations
  - Using ADAM optimiser
- Conditions (choices) guide logical rather than physical flow
  - All expressions evaluated
  - Graph is the same for every event

- Graph implemented with TensorFlow:
  - Supports processing on GPU
  - Wrapped in Keras model:
    - Portable (saving to/loading from disk)
    - Straight forward integration
- Pre-processing:
  - Structure of graph must be static
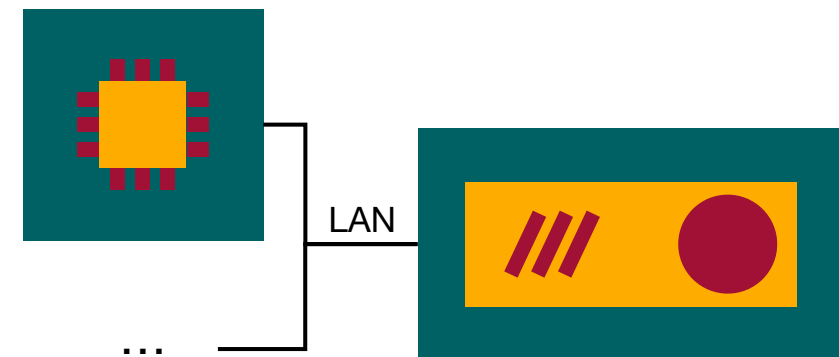  - Requires batching of events in fixed size chunks

## Kinematics

trino momentum:

and W mass

on kinematics:

lex (fit) solution

- Testing two different hardware s
  - Pre-processing always on work

**Experimen...**

- Two different

**On-board**:
- One GPU on each processing n

**Legend:**

| Array | Analytical Operator |
|-------|---------------------|
| Tensor | Fitting Operator |

COM...

Optim... Place... $p_x^2$ $p_y^2$ Ch... pz

## Set... On-board

**Server:**

ware setups:

each de

PU Server

many sing

Connected

rver

- All computations on one computing node

odes a LAN

- Cluster scenarios:
  - CPU-only setup

PU:

ffset for loading and building the graph

rent setups: Each worker has own GPU

during computation

p before computation

ations:

eon Silver 4216

- Pre-heating the G time wo di Setu

**Warm**: Set

Technical specific
- CPU: 2x Intel X
- GPU: NVIDIA (
- Network (LAN)

**Res** Server

- For a typical analysis with O(100M) events:
  - Parallel computing graph >100x faster than typical exec
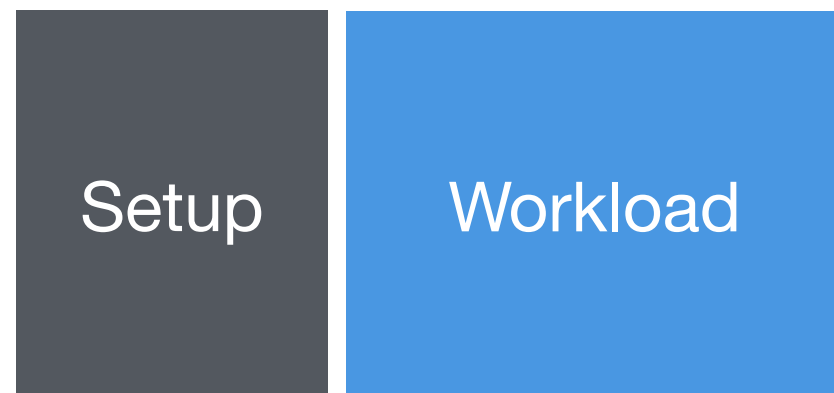  - Speed: On-board GPU > Server GPU = On-board CPU

Graph evaluations on central CPU

Accessed over network (1Gbit/s)

Using TensorFlow model server

Cluster scenario:
- Not each worker has own GPU

Runtime / Event [s]

$10^0$
$10^{-1}$
$10^{-2}$
$10^{-3}$
$10^{-4}$
$10^{-5}$
$10^{-6}$

Legend:
- Graph: On-board CPU (C...
- Graph: On-board GPU (C...
- Graph: Server GPU (Cold...
- ... Graph Server GPU (Warm...
- --- Typical Executable

• Sp

100

10⁻

Dennis Noll - 14.02.22

- Graph need to be built before evaluation

- Takes constant time $\mathcal{O}(10s)$
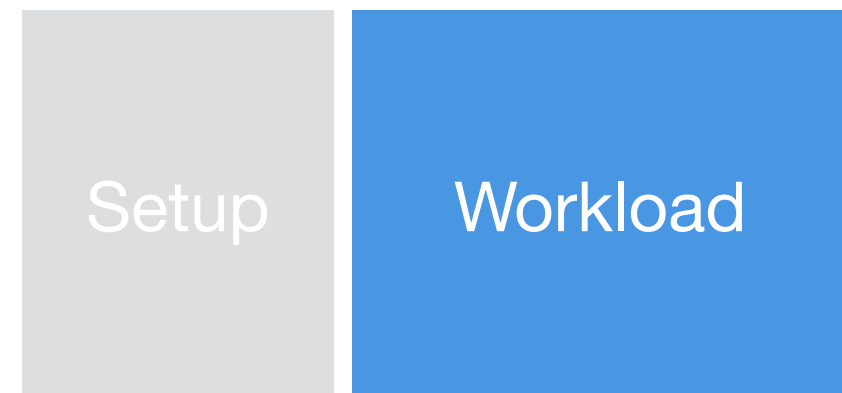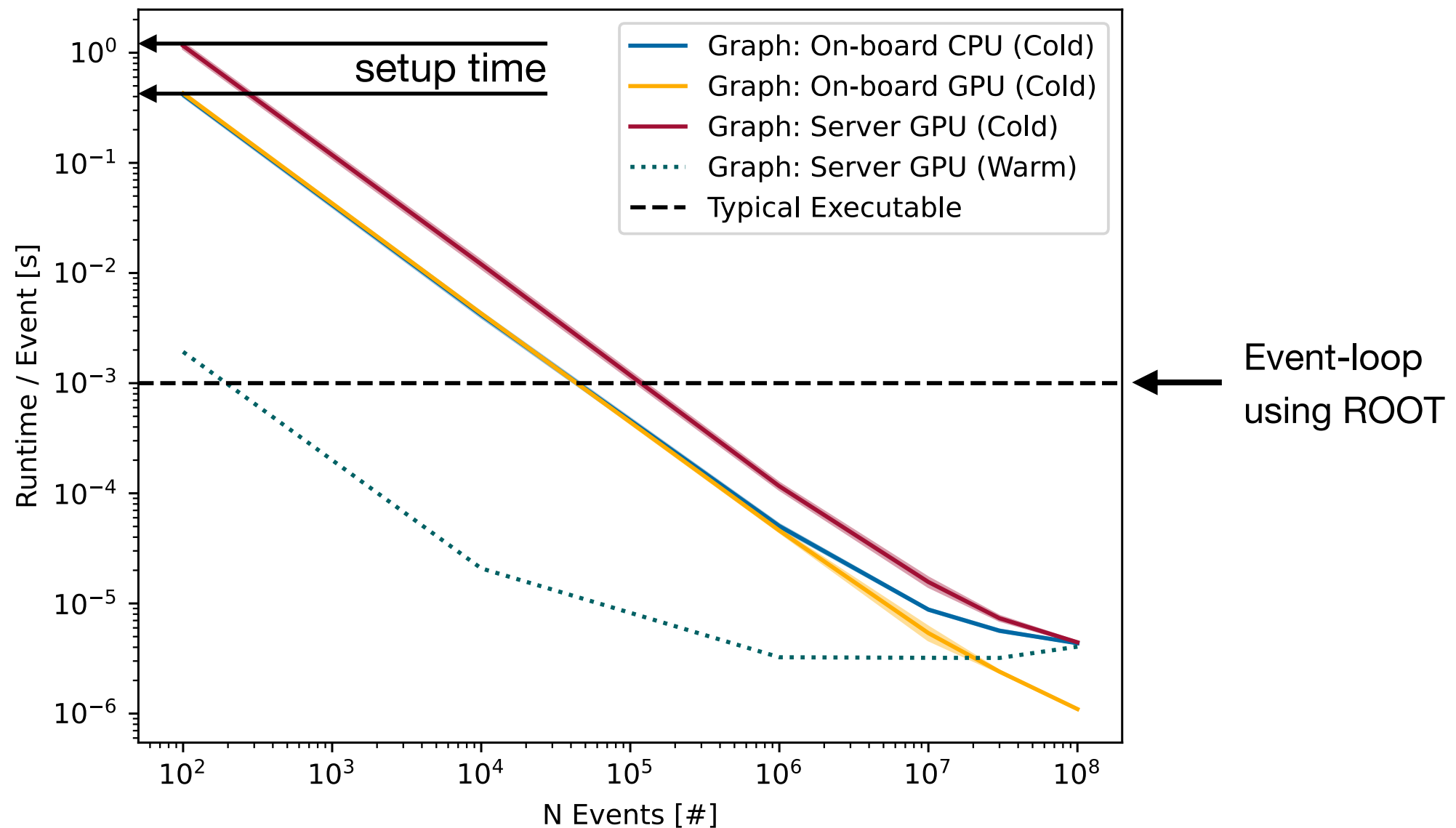
- Testing two scenarios:

**Cold**



Runtime

**Warm**

Runtime

- Setup included in runtime
- Represents:
  - Ad-hoc computation on worker

- Setup not included in runtime
- Represents:
  - Cluster with central GPU

- For typical analysis ($10^8$ events):
  - Graph 100x faster than typical executable
  - On-board GPU: Fastest but also most expensive
  - Server GPU: Saturates due to limited network speed

## ...ations and Graphs

- Future HEP analyses must be fast and resource efficient
- ...g to vectorised computations
  - Use vectorised computations and parallelism

## Computing Graph

- Key functionalities are part of the g...
  - ...cal operations
  - ...ting models
  - ...ser & placeholders for fit...
  - ...operations:
  - ...p with 100 iterations
  - ...aptive optimisation (Adam...
  - ...evaluated during execut...

**Try on** 🌀 **binder**
**https://bndr.it/8yw3z**

**...h:**
...nts

REAL

$p_Z^1$

$p_Z^2$

Choice

Event Kinematics