# Towards GPU-based tracking in ACTS and ATLAS

Stephen Nicholas Swatman
May 3, 2024
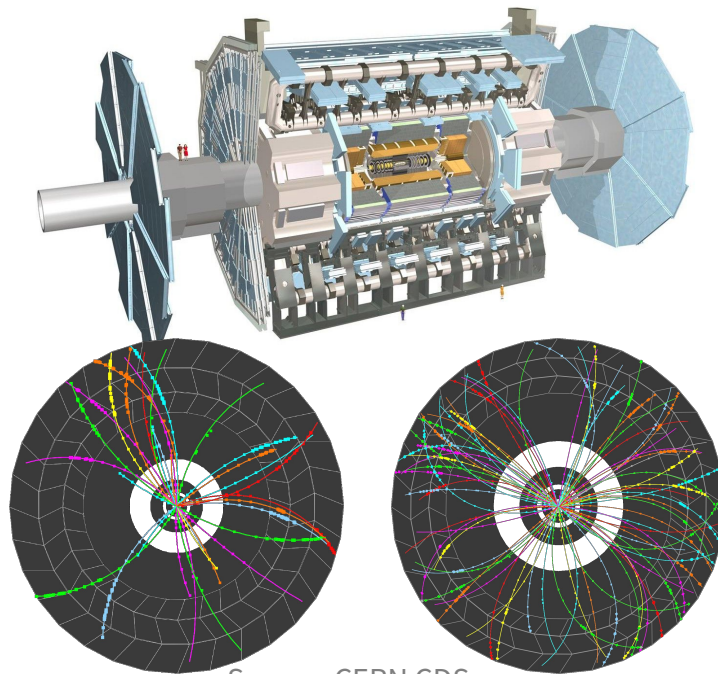
# Introduction – About Me & My PhD

- Finished a MSc at **UvA** + **VU** in 2019, graduation project at **Nikhef**
- Started a PhD in **CS** at **UvA (PCS)** + **CERN** in 2020
- Now writing up – aim to defend in **autumn**
- *Golden staple*-like PhD thesis
- In this talk: very brief overview of my work, and some insights into **GPU tracking**
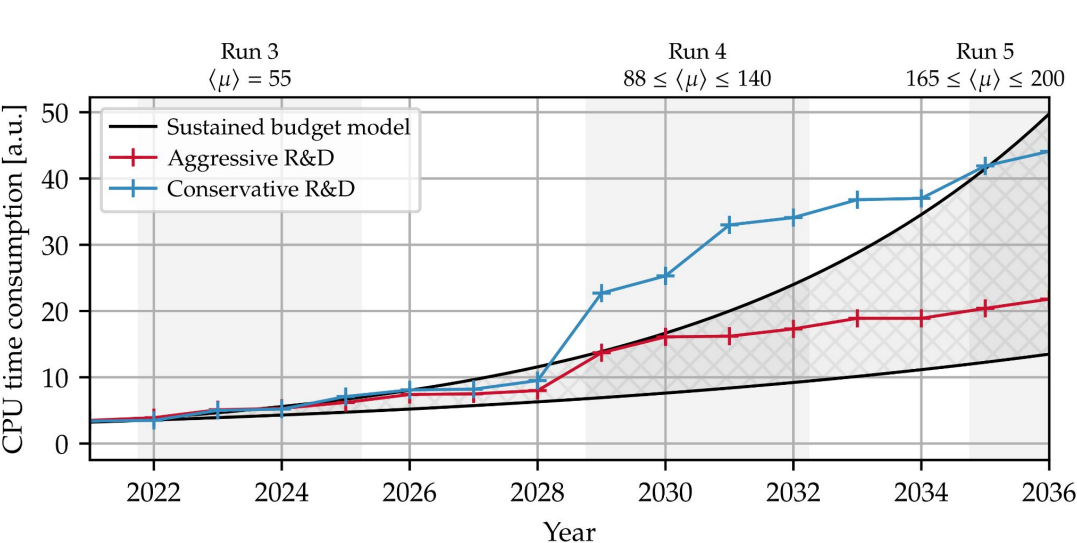
# Problem Statement – Track Reconstruction

- Increased pile-up threatens to make track reconstruction infeasible for HL-HLC, FCC
- More **efficient compute** → more **available resources** → more interesting **physics**
- Compute complexity scales ~$O(\mu^2)$
- This will require research into novel **algorithms**, novel **hardware**, etc.
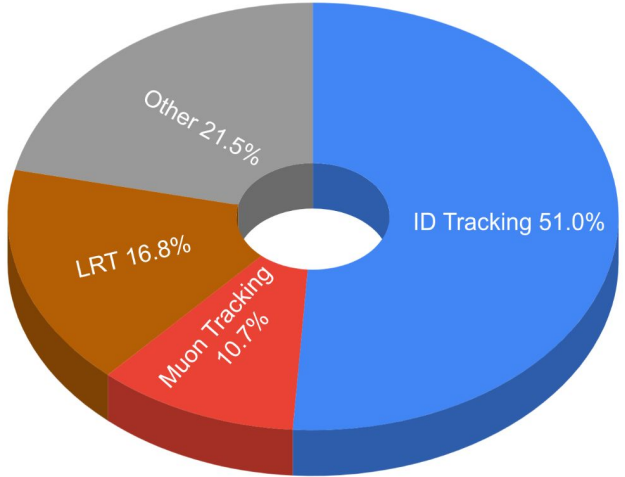- In my case: how can this run on **GPUs**?



Source: CERN CDS

# Problem Statement – Track Reconstruction



Numbers by ATLAS, plot by me



10.1007/s41781-023-00111-y

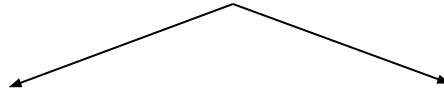# Massively Parallel Computing

**€10,000**

AMD EPYC 9554
64 cores
360W TDP

NVIDIA RTX 6000
18,176 cores
300W TDP

# Massively Parallel Computing



**CPU-like μarch**

**GPU-like μarch**

# Massively Parallel Computing

Can be programmed
independently

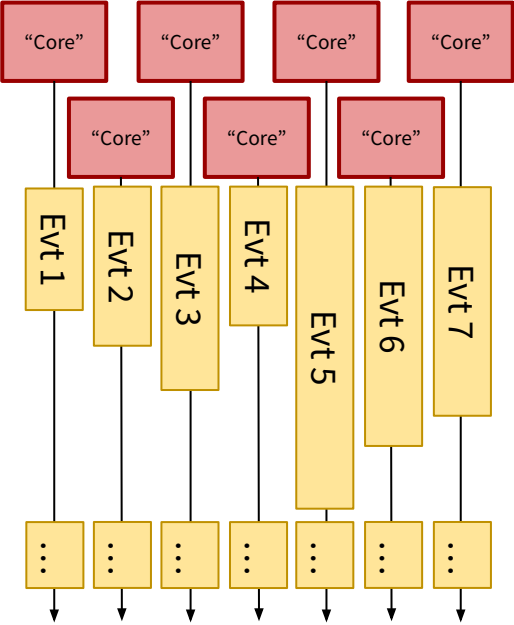| | Cores | × | Cycles/s | × | FLOP/cycle | = | FLOP/s |
|---|---|---|---|---|---|---|---|
| AMD EPYC | 64 | | 3.75B | | 48 | | 11.5 TFLOP/s |
| | | | | | | | ↓ **8×** |
| | 18,176 | | 2.51B | | 2 | | 91.1 TFLOP/s |

Must be carefully
programmed in lockstep

# Massively Parallel Computing
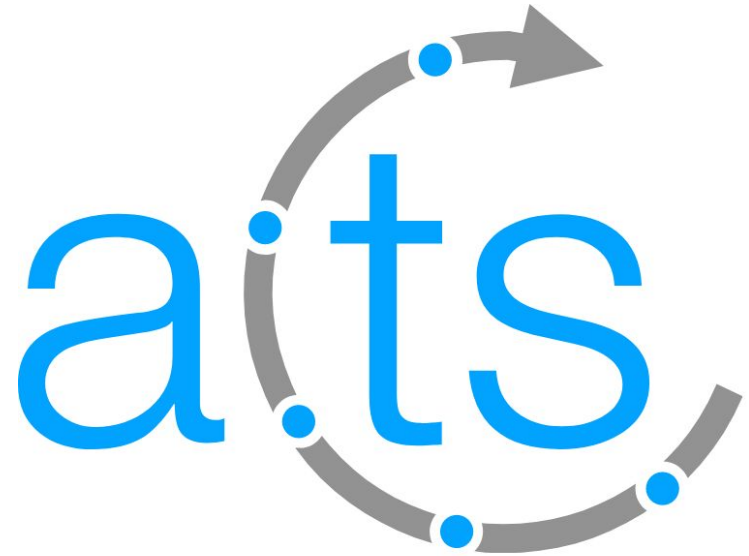


**CPU-like μarch**

**GPU-like μarch**

Can track reconstruction be implemented efficiently on massively parallel systems?

# ACTS – A Common Tracking Software
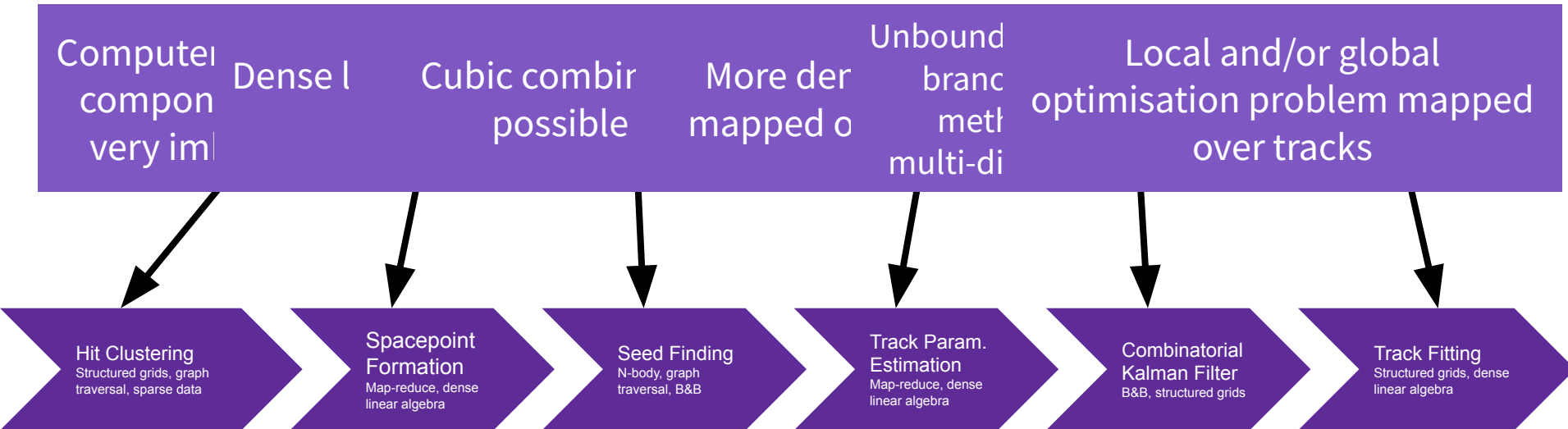
- **ACTS** is **A Common Tracking Software**
- Goals are to be…
    - **Feature complete**
    - **Well-written**
    - **Extensible**
    - **Experiment-agnostic** (sPHENIX, ePIC, etc.)
    - an **R&D platform**
- Originally based on **ATLAS tracking**
- Now being migrated **back into Athena**

What are the challenges in developing track reconstruction algorithms to massively parallel architectures?
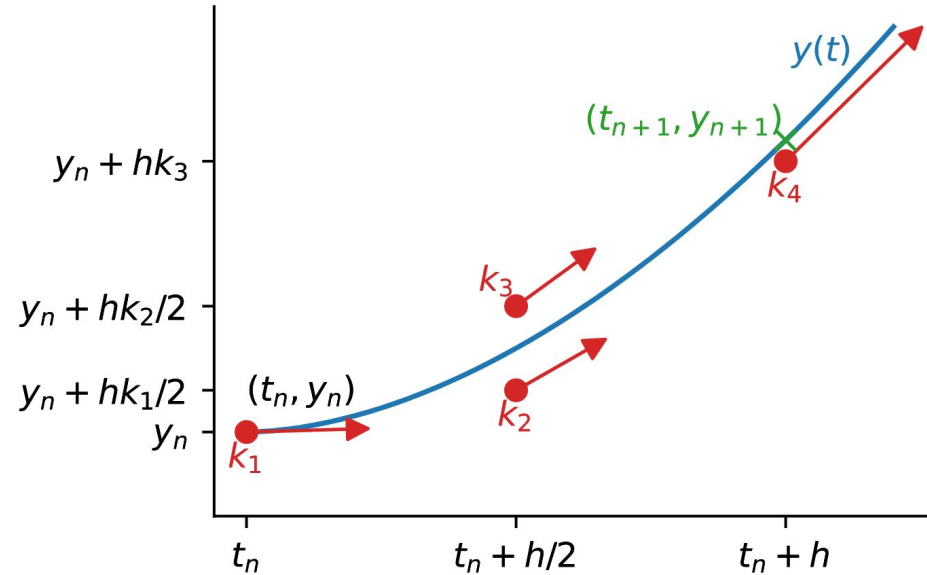
# Problem Statement – State-of-the-Art

- Track reconstruction consists of a **task graph** of structurally different algorithms
- Here denoted according to their "13 dwarves" classification
- Varying degrees of **complexity** and **challenge**…

Computer
compon
very im

Dense l

Cubic combin
possible

More der
mapped o

Unbound
branc
meth
multi-di

Local and/or global
optimisation problem mapped
over tracks

**Hit Clustering**
Structured grids, graph
traversal, sparse data

**Spacepoint Formation**
Map-reduce, dense
linear algebra

**Seed Finding**
N-body, graph
traversal, B&B

**Track Param. Estimation**
Map-reduce, dense
linear algebra

**Combinatorial Kalman Filter**
B&B, structured grids

**Track Fitting**
Structured grids, dense
linear algebra

How can structured grid data be represented in order to maximize the efficiency of arbitrary computations?

# Research – Vector Fields – Amuse Bouche

- Most experiment **magnetic fields** are fairly homogeneous, but not quite
- Those irregularities matter when e.g. **propagating particle motion**
- Storing ATLAS B-field: **~200 MB**
- **Accessed** millions of times per second!
- How do we do that **quickly**?
- We don't even know how to do than on CPUs, and if we did it would not translate!



Source: ACTS Project

# Research – Vector Fields

- **Multi-dimensional data** is everywhere HPC
  - Magnetic fields
  - But also CFD, lattice QCD, etc.
- Must be accessed very frequently in **hard-to-predict patterns** for iterative numerical methods
- Increasing **cache efficiency** can improve **performance** for these kernels
- **Design space is large**
  - Many **functional** and **non-functional** properties come into play

Source: Moritz Lehmann

# Research – Vector Fields

- Interpolation methods (F & EF)
  - NN, linear, cubic, etc.
- Boundary checking (F & EF)
  - Wrap, mirror, zero, etc.
- Array layout (EF)
  - Row-major, column-major, *Morton*, etc.
- Coordinate transformation (F & EF)
  - Affine, polar, etc.
- Storage location (EF)
  - Main memory, CUDA memory, texture memory, etc.
- Potentially many more!
- Across many **devices** and with many access patterns



1D nearest-neighbour

Linear

Cubic

2D nearest-neighbour

Bilinear

Bicubic

Source: Wikipedia

# Research – Vector Fields

- We developed a **category-theoretical** method for **decomposing** this design space (like on the previous slide)
- Can be re-composed at **compile time** with zero run-time overhead
- Allows our method to serve as both a **benchmarking suite** and **design space exploration tool**…
- …as well as a **library** for implementing heterogeneous multi-dimensional arrays

ATLAS magnetic field

$\mathbb{R}^3$ $\mathbb{R}^3$

# Research – Vector Fields

- This work was published in **ICPE'23**
- And it was nominated as candidate to the **best paper award**!
- https://doi.org/10.1145/3578244.3583723

# Research – Morton Layouts – Amuse Bouche

# Research – Morton Layouts

- The Morton curve provides **balanced locality** compared to row-major and column-major layouts
- Calculated by **interleaving** the bits of the **binary expansions** of the input coordinates!
  - Can be done efficiently on modern commodity hardware
- But what if you were to interleave the bits in **arbitrary patterns**?

$$f(5, 3, 4) = f(101_2, 011_2, 100_2) = 1010101110_2 = 342_{10}$$

# Research – Morton Layouts

$$f(1011_2, 0101_2) = \underline{\vee \begin{matrix} 01000101_2 \\ 00100010_2 \end{matrix}} = 103_{10}$$
$$01100111_2$$

$$f(1011_2, 0101_2) = \underline{\vee \begin{matrix} 01000110_2 \\ 00010001_2 \end{matrix}} = 87_{10}$$
$$01010111_2$$

# Research – Morton Layouts

- Turns out this gives you *very* large families of array layouts, all of which have **different cache properties**!
- We propose that **evolutionary algorithms** can be used to efficiently explore this design space
- We show that we can significantly **improve performance** – up to 10 times in extreme cases
- Automated, **problem-agnostic** optimisation method!



(a) [0,0,0,1,1]   (b) [0,0,1,0,1]   (c) [0,0,1,1,0,1]   (d) [0,0,1,1,1,0]

(e) [0,1,0,0,1,1]   (f) [0,1,0,1,0,1]   (g) [0,1,0,1,1,0]   (h) [0,1,1,0,0,1]

(i) [0,1,1,0,1,0]   (j) [0,1,1,1,0,0]   (k) [1,0,0,0,1,1]   (l) [1,0,0,1,0,1]

(m) [1,0,0,1,1,0]   (n) [1,0,1,0,0,1]   (o) [1,0,1,0,1,0]   (p) [1,0,1,1,0,0]

(q) [1,1,0,0,0,1]   (r) [1,1,0,0,1,0]   (s) [1,1,0,1,0,0]   (t) [1,1,1,0,0,0]

# Research – Morton Layouts

- This was accepted to **ICPE'24**
- To be presented in **London** next week!
- https://doi.org/10.1145/3629526.3645034 (not yet active)

How can the effects of thread imbalance in SIMT workloads be modelled and how can they be mitigated?
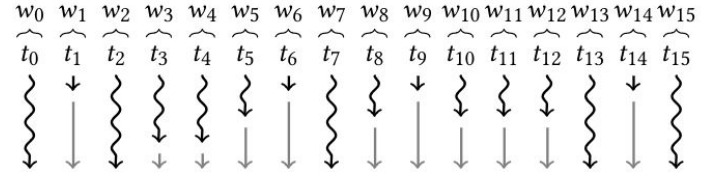
# Research – Thread Imbalance – Amuse Bouche

- We know we cannot assign one event to one thread: **lockstep execution**
- Turns out we also cannot assign one module to one thread: **too imbalanced**
- We will need **fundamentally** different algorithms for **clustering**
- But can we **predict** what will work?
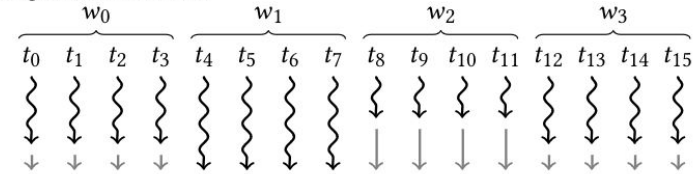  - And can we use the graph on the right?
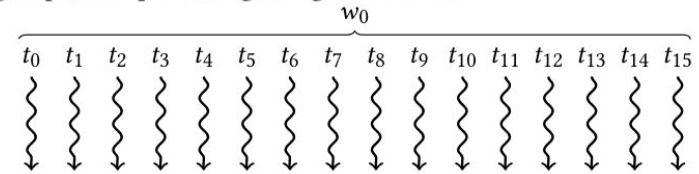
# Research – Thread Imbalance

- A lot of our kernels are **imbalanced**: different threads execute **different amounts of work**
  - Different cluster sizes
  - Different branching factors
- This is not a problem on CPUs, but **strongly impacts GPU performance**
- Can be mitigated using **thread coarsening** or **thread refinement**
- But what is the **performance impact** of these techniques?



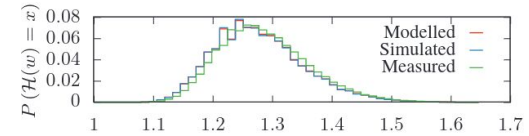(a) Per-thread granularity: each thread in a thread group processes a single unit of work.

(b) Intermediate granularity: threads in a thread group form sub-groups, each processing a single unit of work.

(c) Per-group granularity: all threads in a thread group process a single unit of work.
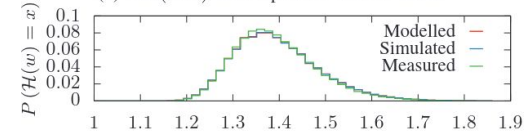
# Research – Thread Imbalance

- We propose a **statistical model** that can infer the overhead of SIMT execution
- Metric for "*how suitable for GPU execution is this workload*"
- Uses **only** distribution of **thread load**: no hardware details, etc.
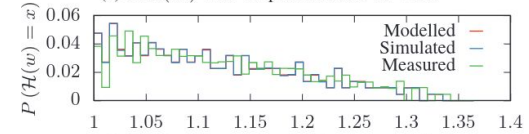- Allows **early evaluation** of feasibility of different **parallelism strategies**!
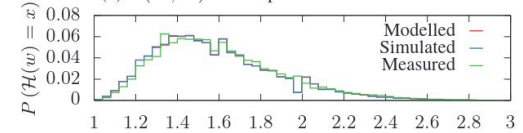


(a) B(40, 0.5) with 16 parallel units of work.

(b) Geo(0.05) with 8 parallel units of work.

(c) Pois(30) with 32 parallel units of work.

(d) U(20, 40) with 2 parallel units of work.

(e) NB(5, 0.3) with 4 parallel units of work.

# Research – Thread Imbalance



- This work was published in **MASCOTS'22**
- Which was nice and also in Nice
- https://doi.org/10.1109/MASCOTS56607.2022.00026

How do the extra-functional properties of novel track reconstruction algorithms compare to state-of-the-art solutions?

# traccc – GPU Tracking Demonstrator

- Now: back to the application under study
- We have developed a **track reconstruction chain** for massively parallel devices
- Integrates **novel algorithms**
- Tracking on **TrackML** and **ODD**-like detectors

# traccc – GPU Tracking Demonstrator

# ACTS Project – Subprojects

- R&D consists of many **subprojects** for **HEP** and **HPC** in general
- **traccc**: tracking demonstrator
- **algebra-plugins**: linear algebra
- **detray**: detector description
- **vecmem**: memory management
- **covfie**: vector field storage
  - Plasma physics at HZDR

# Key Point – Reproducibility

- Aim: to write software and develop methods that can be **used and improved on**
  - LHC lifetime: 17 more years
- Try to avoid "PhDware": software that becomes unusable after the end of the PhD
- **Artifact** evaluation tracks take extra effort but reward handsomely
  - Nice **stickers**
  - In some cases taken into account for reviews, rebuttals, etc.
- Also served on AE committees for SC, CGO, ICPE, and ICPP

Source: ACM

# Results



Maximum throughput

- *Preliminary* results on **TrackML** data show that our GPU-based solutions work well
- **Outperform** similarly priced CPUs at **higher pile-up values**
- **Factor ~10** gain in throughput

Source: Guilherme Metelo Rita de Almeida

# Outstanding Challenges

- So far tested only on **simple geometries**: how do we integrate e.g. **ITk**

- **Combinatorial Kálmán Filter**: important step with high combinatorics

  - How do we **distribute branches** over threads?

- Integration into **Athena** is ongoing work

- How to **schedule** and **place** algorithms?

  - GPUs have **separate memories**: transfers are not free

# Research – Throughput Models

- Our scheduling is an open problem, but can we somehow estimate the **throughput** of our **task graph on heterogeneous systems**…
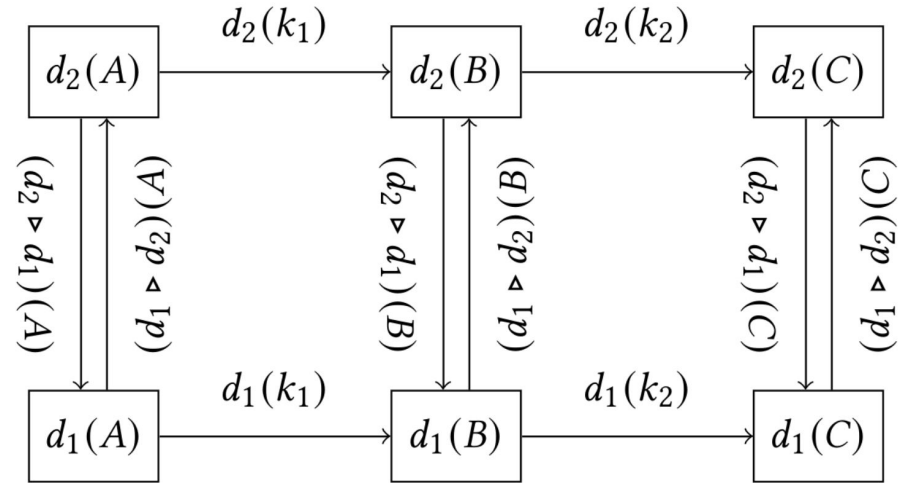- …using only the throughputs of the **individual kernels**?
- We propose that we can create an **optimistic upper bound** for this based on work in the data flow community
- Using **linear programming** we can solve a resource-constrained **maximum flow problem**!



$$\text{maximise} \quad \sum_{e \in E^-(t)} \boldsymbol{x}_e f(e)$$

$$\text{subject to} \quad \forall v \in T' \setminus \{s, t\} : \sum_{e \in E^+(v)} \boldsymbol{x}_e f(e) = \sum_{e \in E^-(v)} \boldsymbol{x}_e f(e)$$

$$\forall e \in E : 0 \le \boldsymbol{x}_e \le 1$$

$$\forall r \in D \cup I : 0 \le \sum_{e \in Q(r)} \boldsymbol{x}_e \le 1$$

# Conclusions

- ATLAS needs **aggressive R&D** to tackle **high-μ compute challenges**

- **Massively parallel** track reconstruction under development in **ACTS**

- Despite **irregular workloads**, we can exploit **GPUs** well in **TrackML** and **ODD**

- **Performance** is very **competitive** at $\mu \geq 100$

- Complex **geometries**, **scheduling**, and **placement** remain open questions

The ACTS project
https://github.com/acts-project/

Bi-weekly R&D meeting
https://indico.cern.ch/category/16958/