

# Tracking Machine Learning Challenge

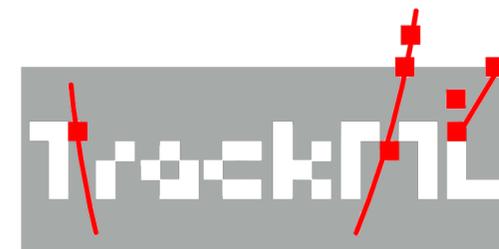
A summary



@SaltyBurger

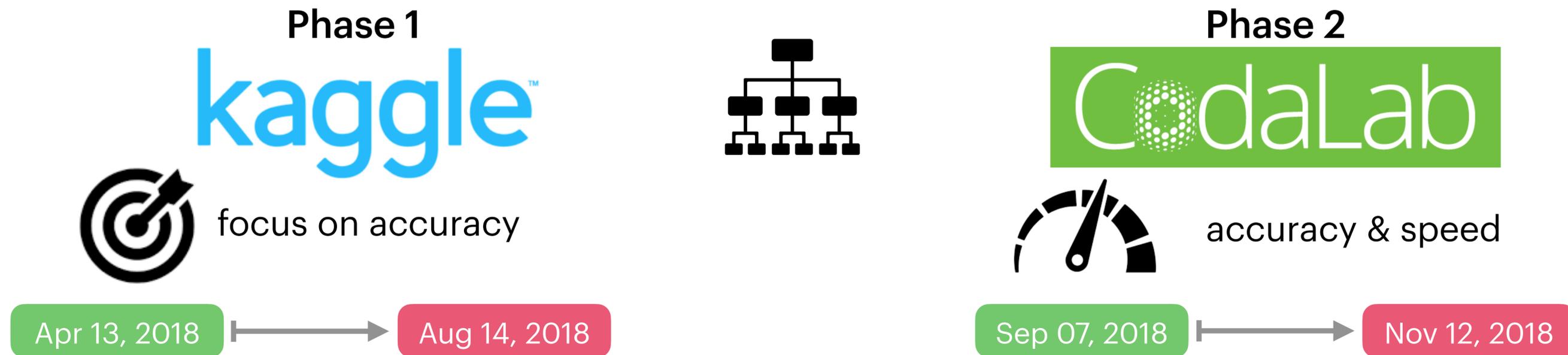


A. Salzburger (CERN) for the



organisers

# Organisation



## Organisation team

Jean-Roch Vlimant (Caltech), Vincenzo Innocente, Andreas Salzburger (CERN), Isabelle Guyon (ChaLearn), Sabrina Amrouche, Tobias Golling, Moritz Kiehn (Geneva University), David Rousseau, Yetkin Yilmaz (LAL-Orsay), Paolo Calafiura, Steven Farrell, Heather Gray (LBNL), Vladimir Vava Gligorov (LPNHE-Paris), Laurent Basara, Cécile Germain, Victor Estrade (LRI-Orsay), Edward Moyse (University of Massachussets), Mikhail Hushchyn, Andrey Ustyuzhanin (Yandex, HSE)



## Summary publications



Particle Tracking Machine Learning Challenge [ [Phase 1](#) ] [ [Phase 2](#) ]

# Sponsors



kaggle



**NVIDIA**

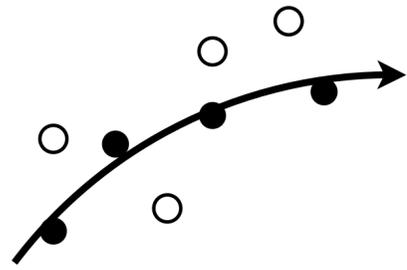


**UNIVERSITÉ  
DE GENÈVE**



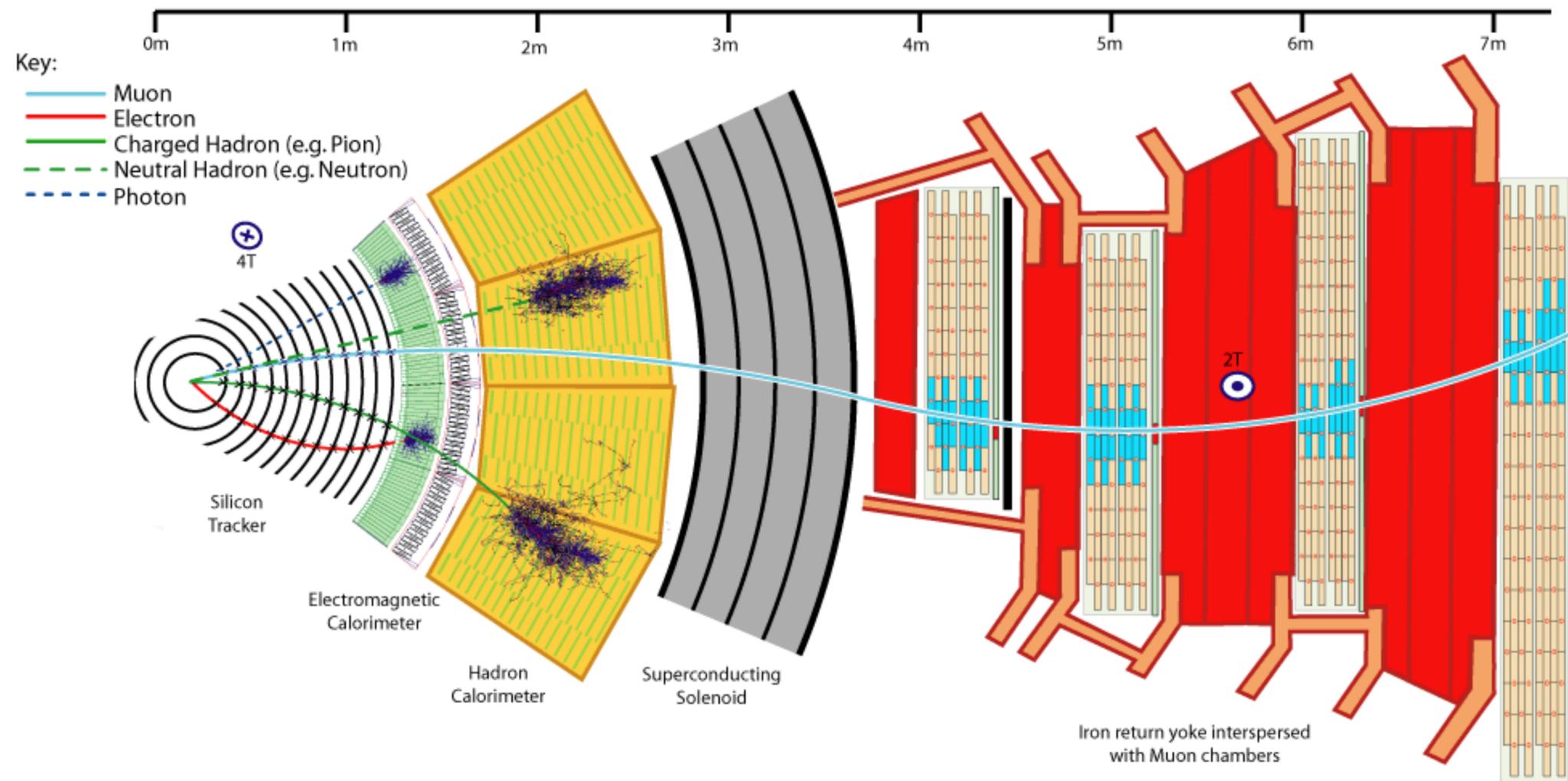
Paris-Saclay  
Center for  
Data Science

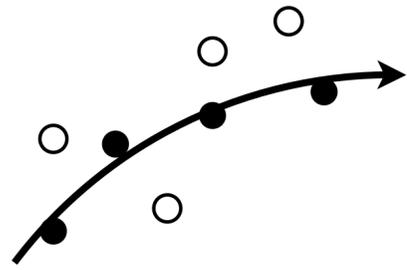




# Particle Tracking

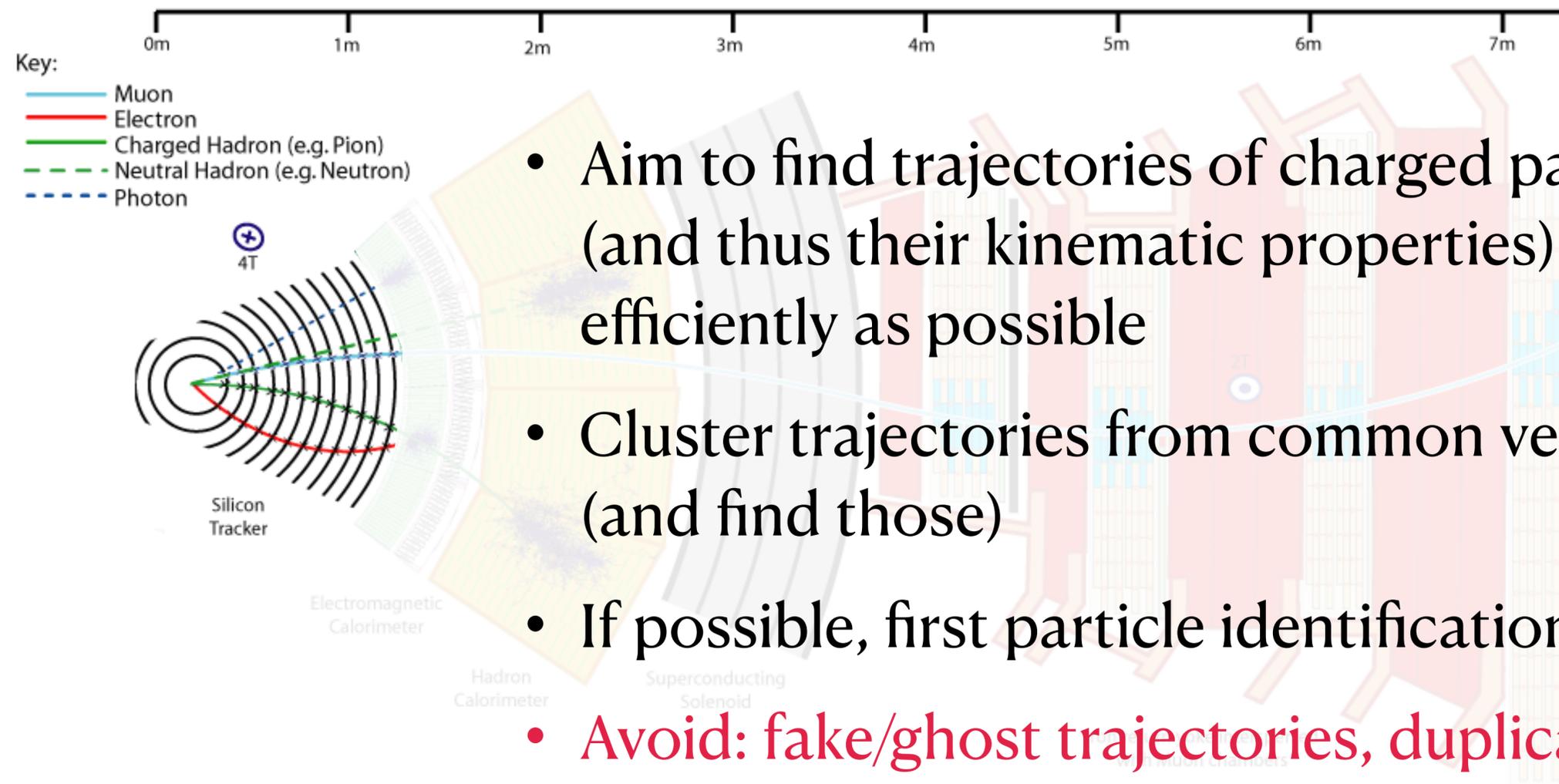
Trajectory and vertex finding in tracking detectors

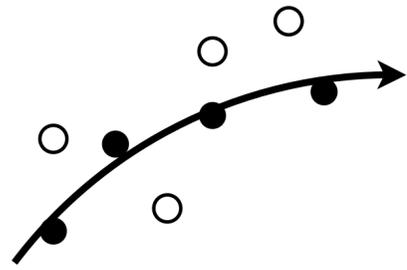




# Particle Tracking

Trajectory and vertex finding in tracking detectors

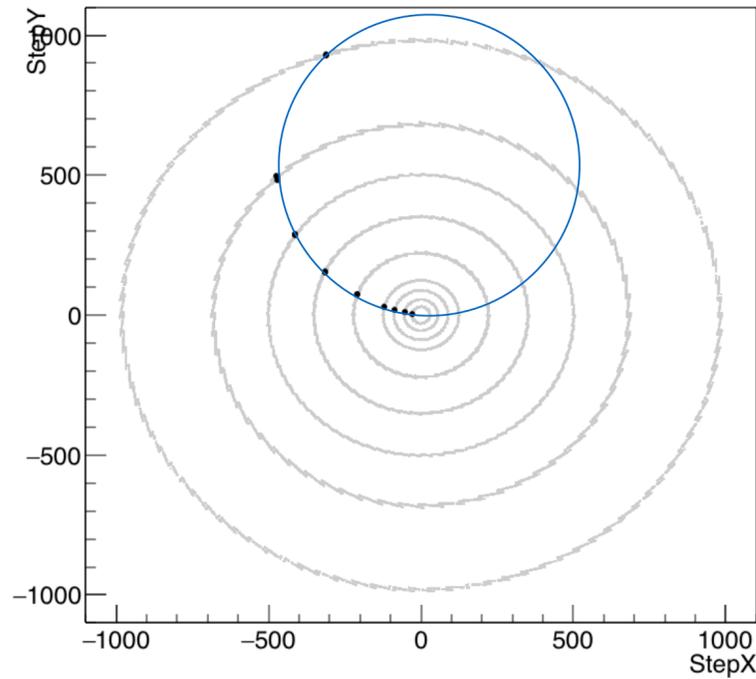




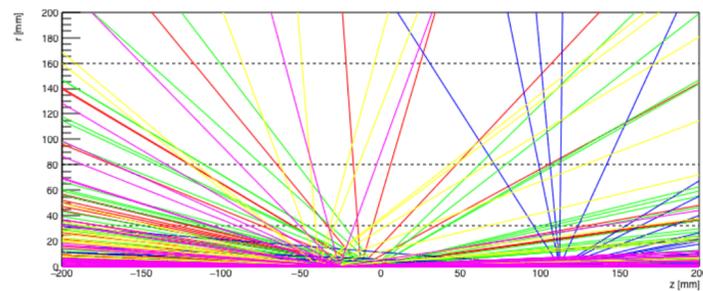
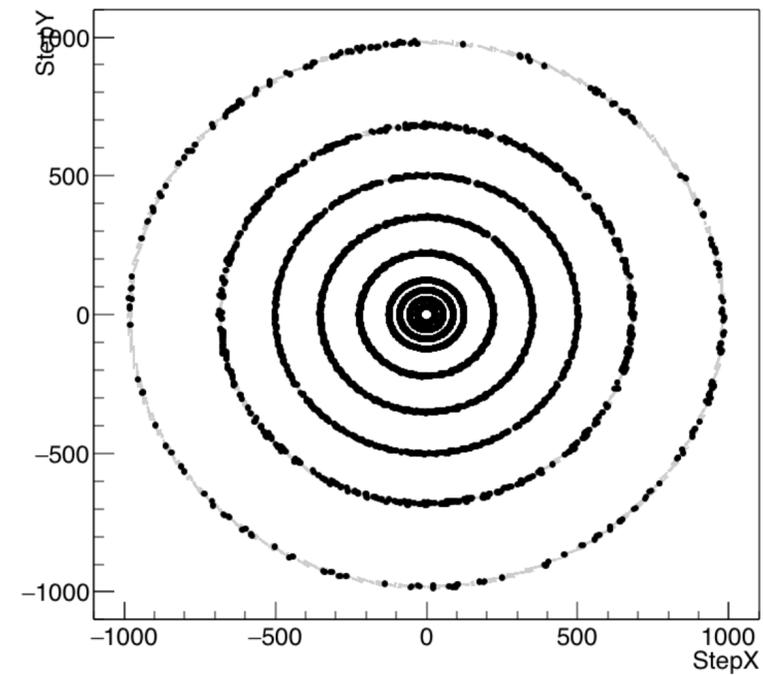
# Particle Tracking

Complexity with increasing pile-up

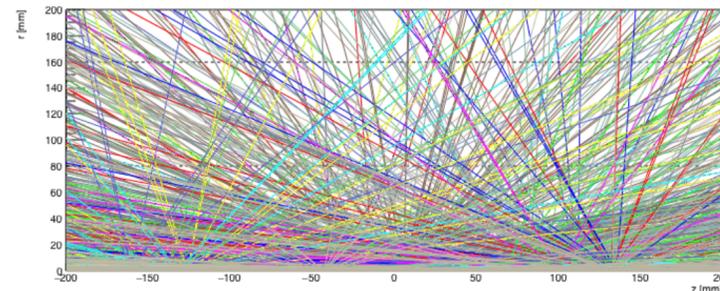
Hits from one particle



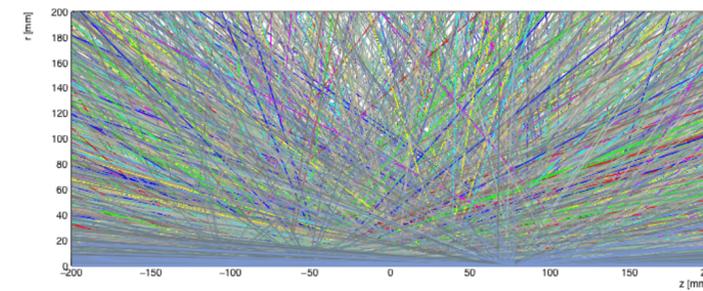
Fraction of hits of a  $\langle \mu \rangle \sim 200$  event



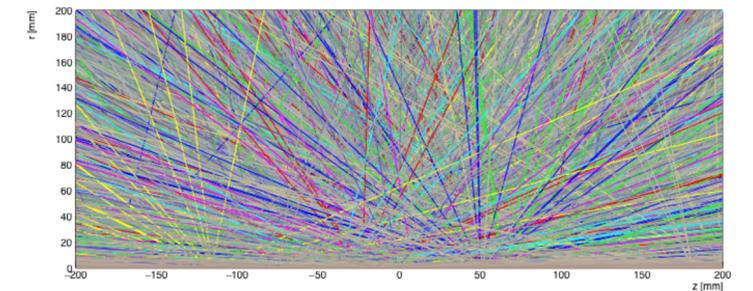
LHC Run-1,  $\langle \mu \rangle \sim 5$



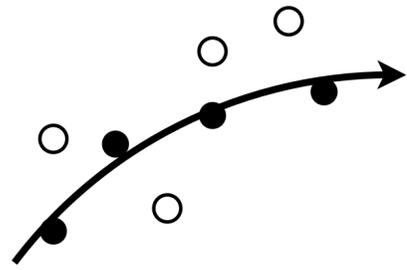
LHC Run-2,  $\langle \mu \rangle \sim 20$



HL-LHC,  $\langle \mu \rangle \sim 200$

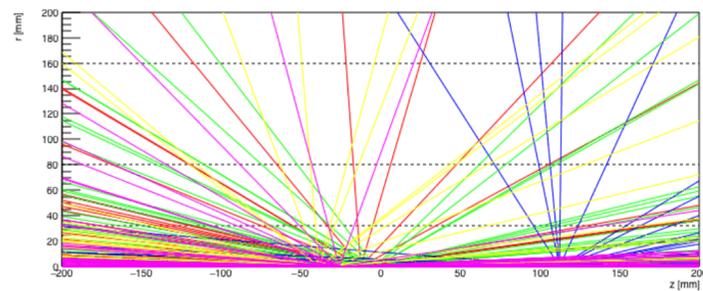
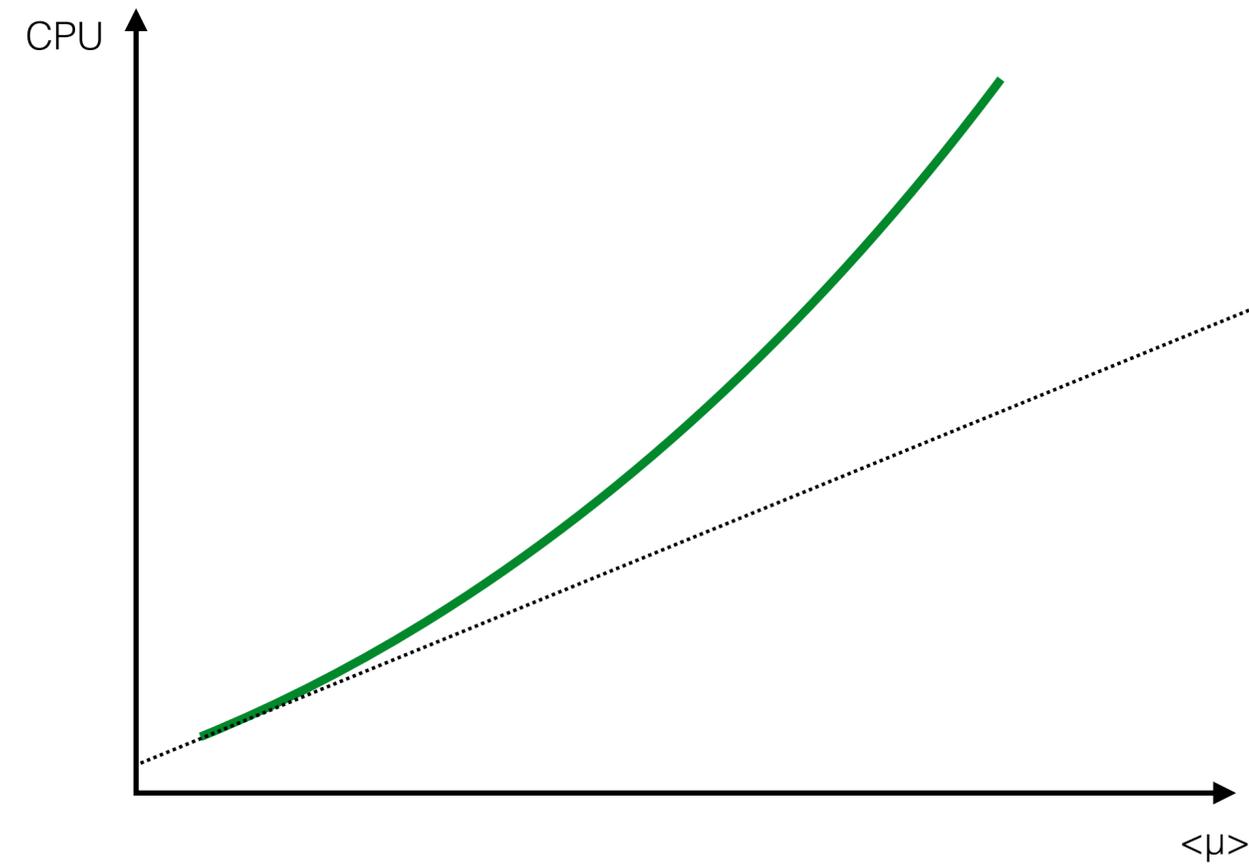


FCC-hh (25ns)  $\langle \mu \rangle \sim 1000$

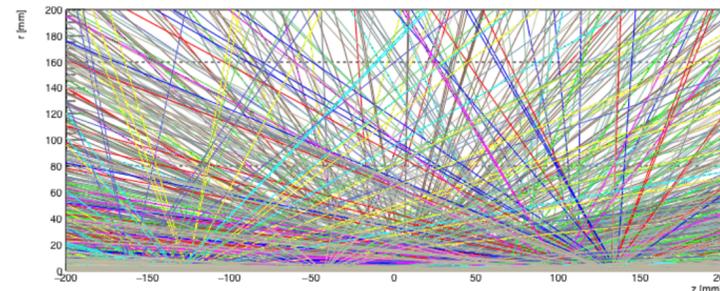


# Particle Tracking

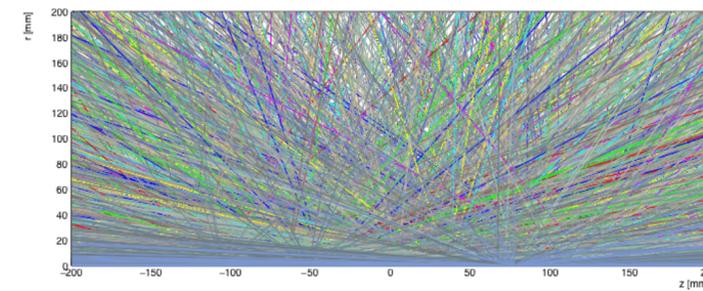
A combinatorial / computational challenge



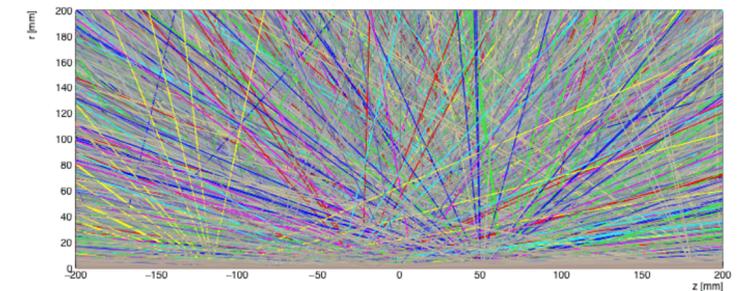
LHC Run-1,  $\langle \mu \rangle \sim 5$



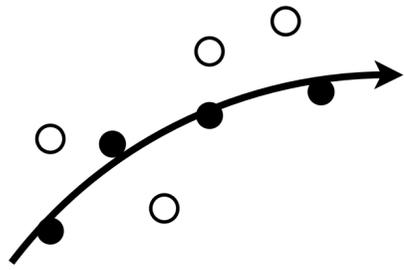
LHC Run-2,  $\langle \mu \rangle \sim 20$



HL-LHC,  $\langle \mu \rangle \sim 200$



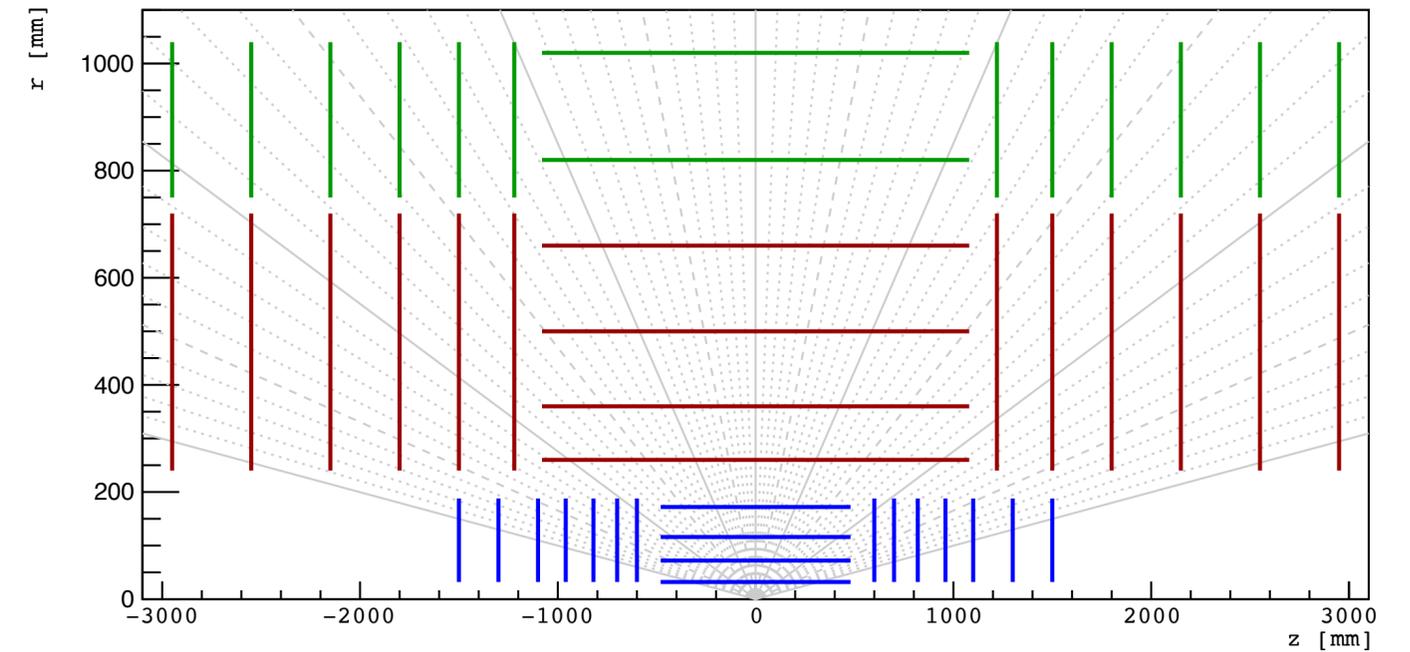
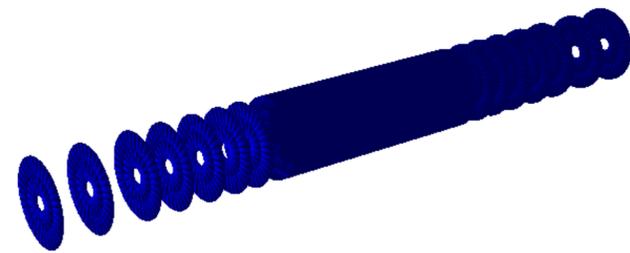
FCC-hh (25ns)  $\langle \mu \rangle \sim 1000$

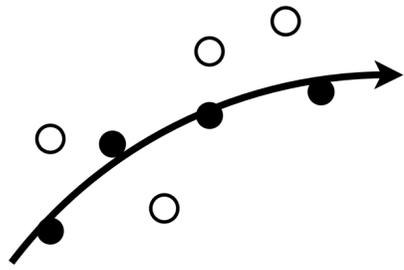


# The Challenge

The TrackML detector - Innermost Pixel System

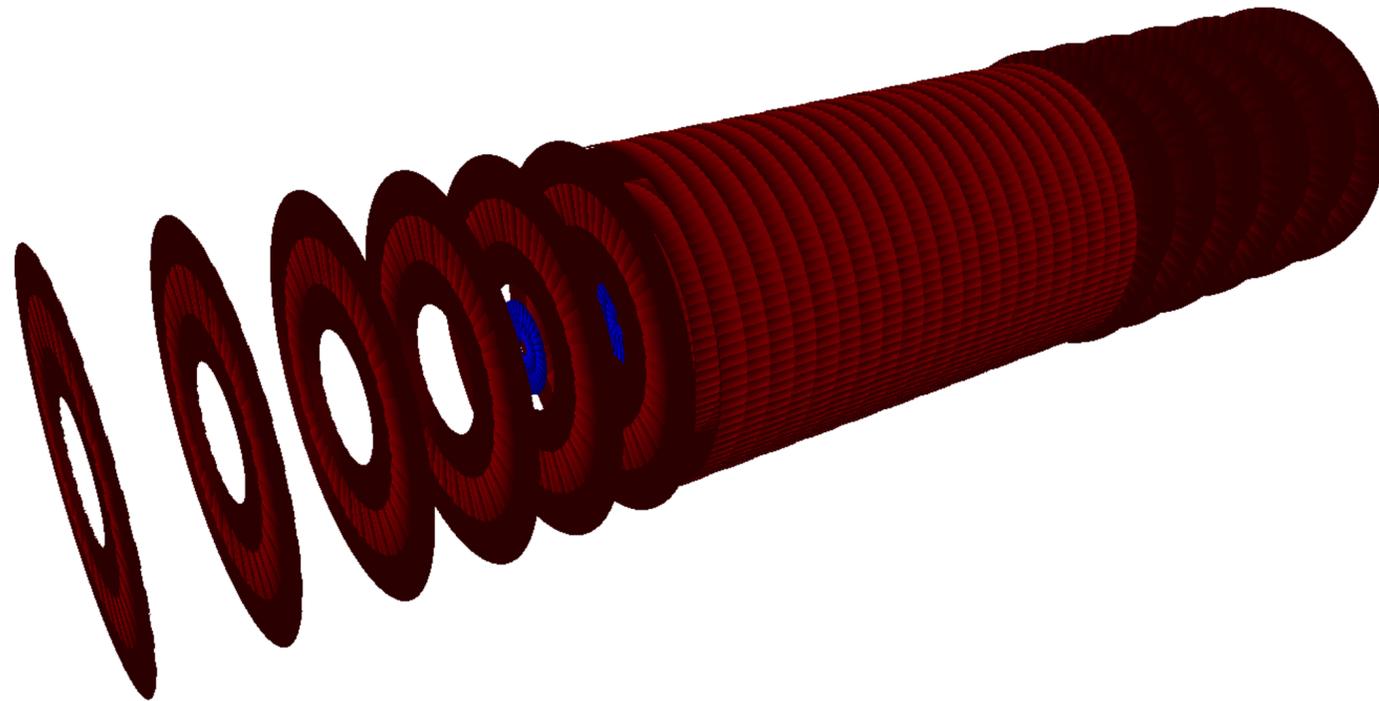
4 barrel-layer central pixel system  
7 EC discs (both sides)



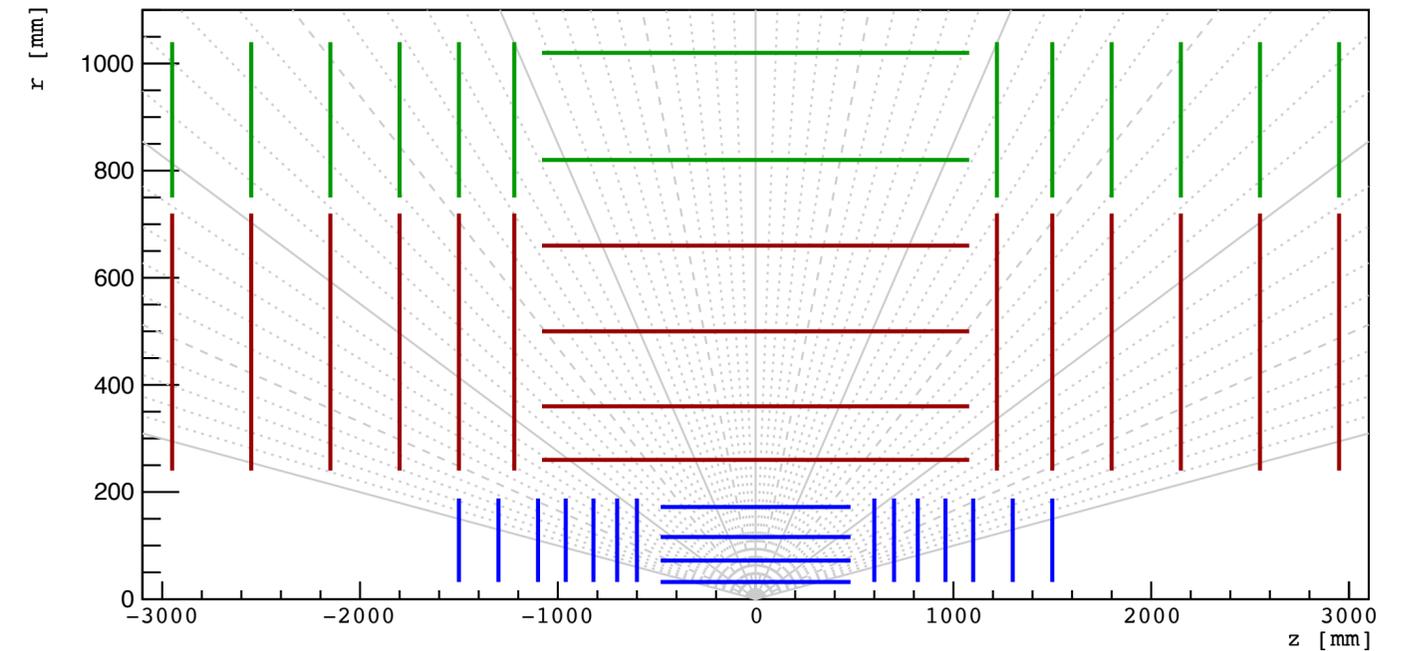


# The Challenge

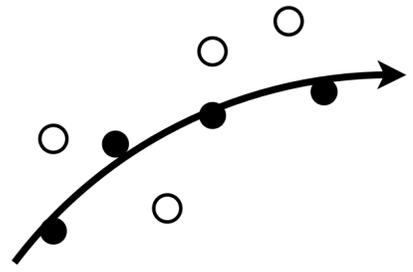
The TrackML detector - Short Strip System



4 barrel-layer central strip system  
6 EC discs (both sides)

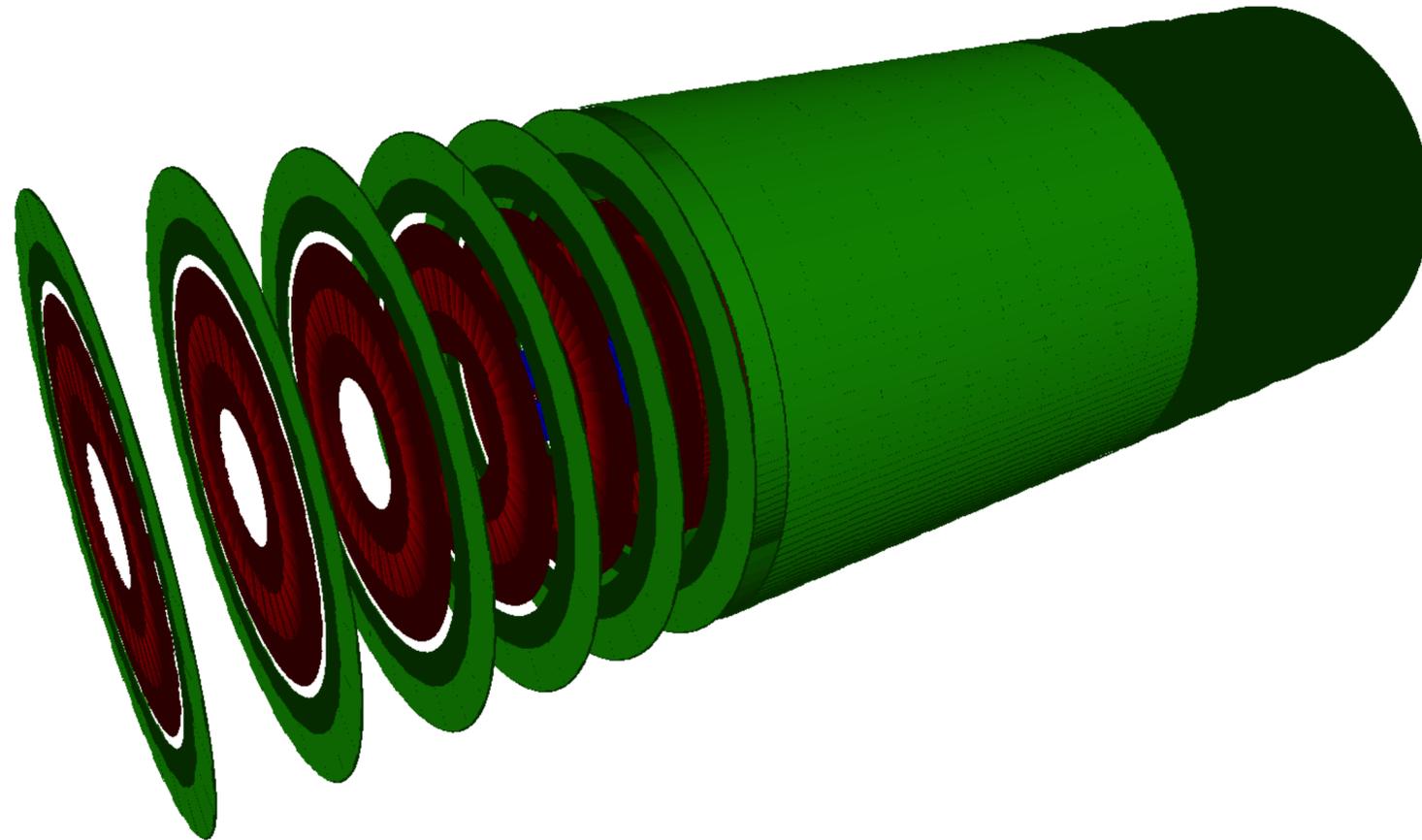


Coverage  $|\eta| < 3$

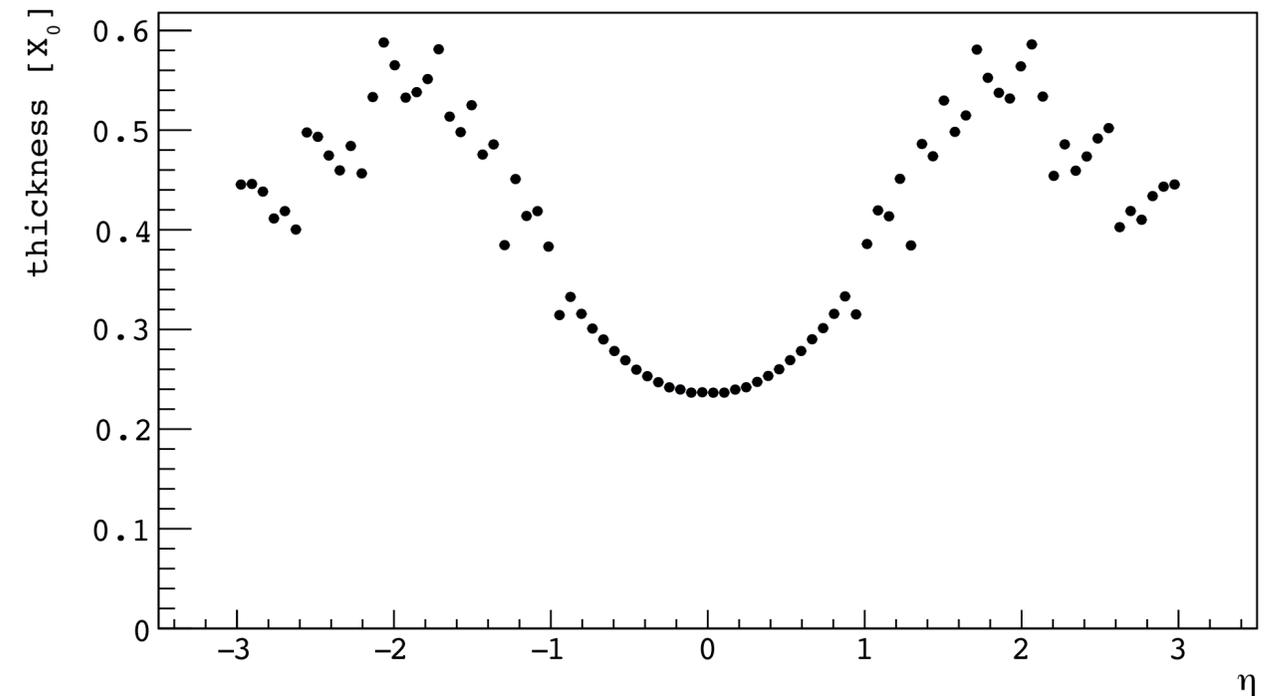


# The Challenge

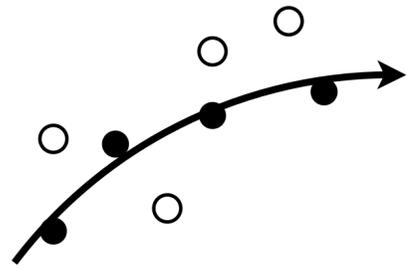
The TrackML detector - Long Strip System



2 barrel-layer outer strip system  
6 EC discs (both sides)



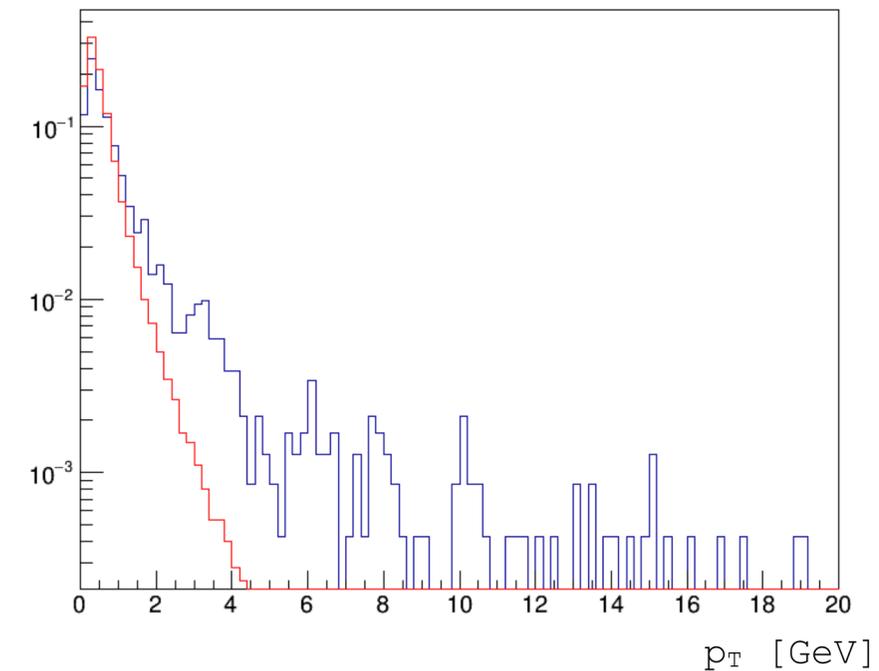
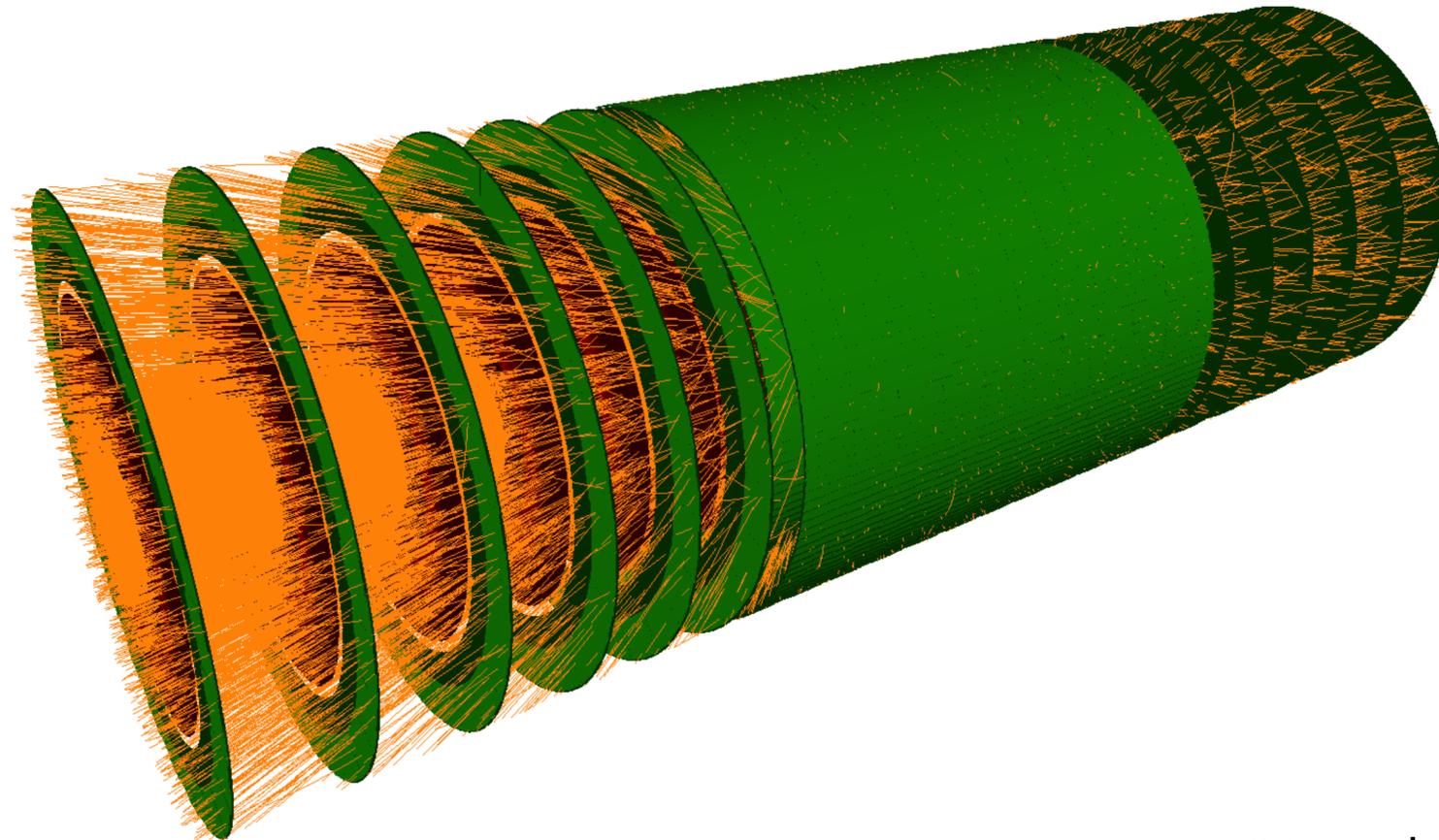
Optimistic material budget



# The Challenge

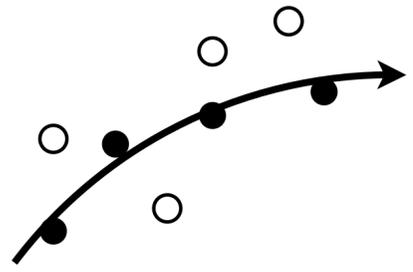
A TrackML event

Signal event: top quark pair production  
Overlaid pile-up:  $\langle \mu \rangle = 200$



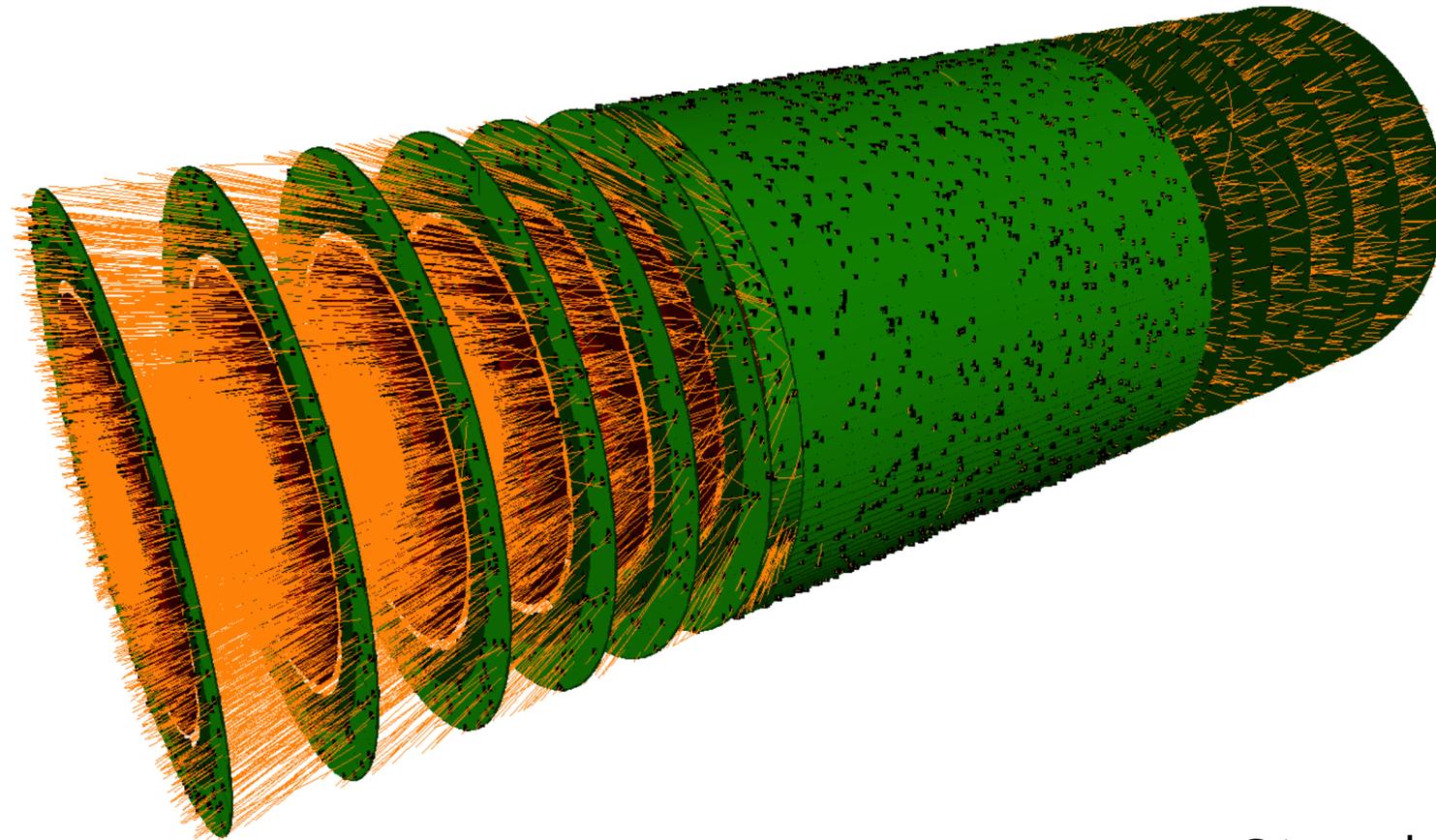
Simulated with [aits](#) fast track simulation

→ 10k particles, 100k hits / event



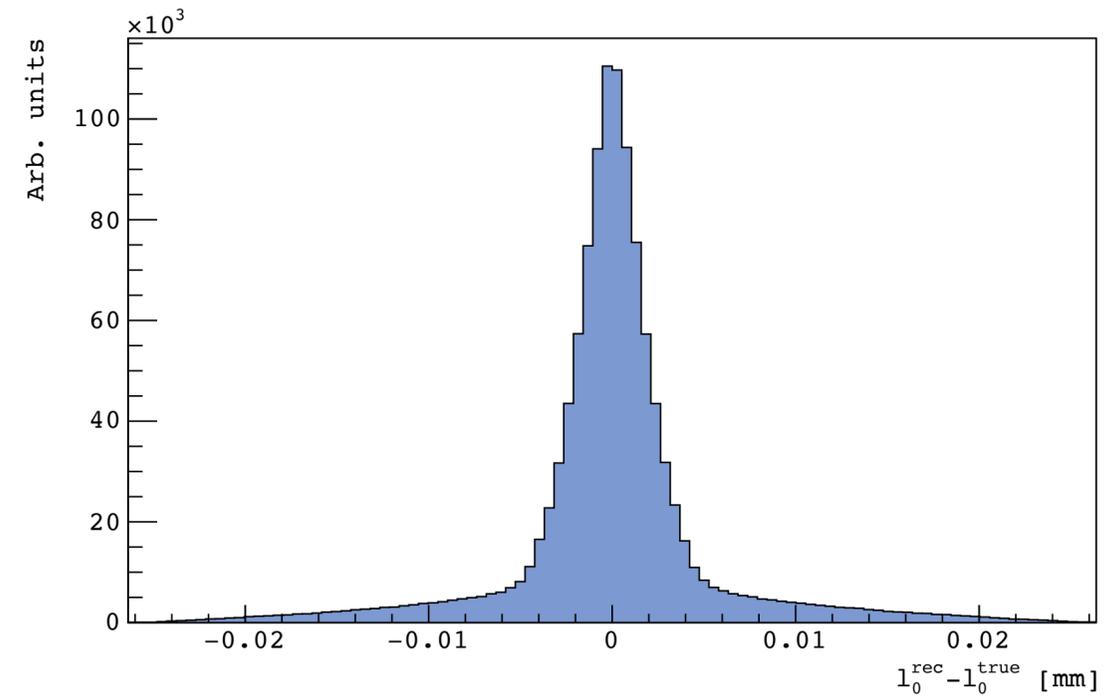
# The Challenge

A TrackML event

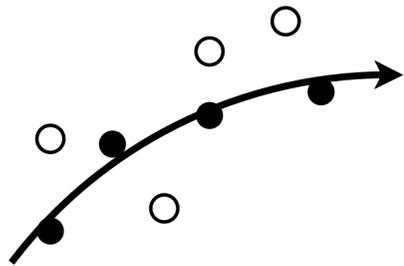


Quasi-realistic resolution

- Using geometric digitisation approach
- Non-gaussian hit residuals

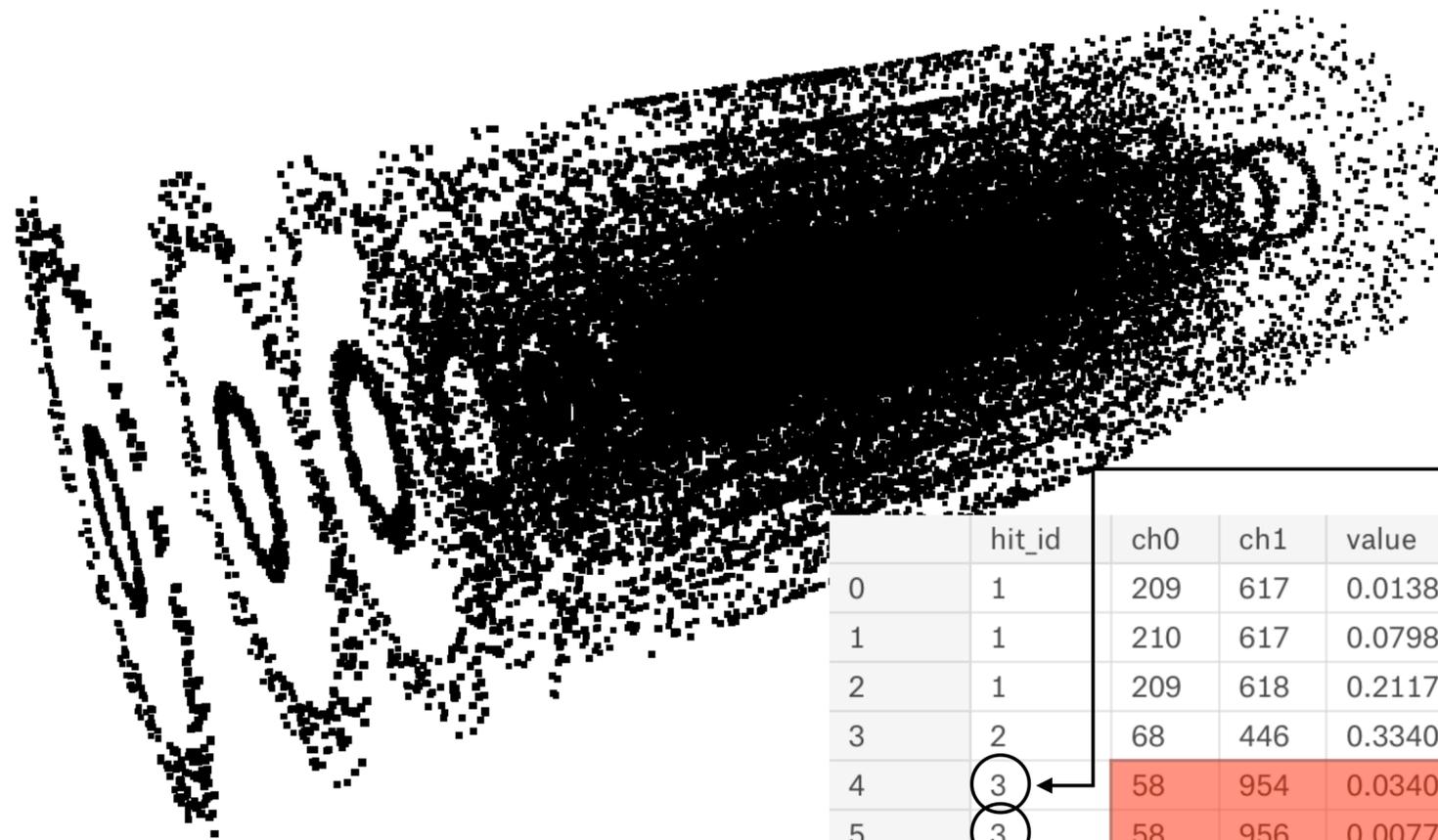


Simulated with [ats](#) fast track simulation



# The Challenge

A TrackML event - The Dataset

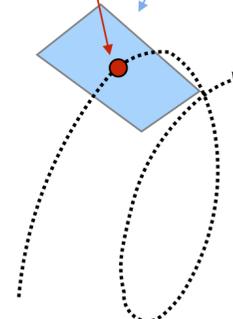
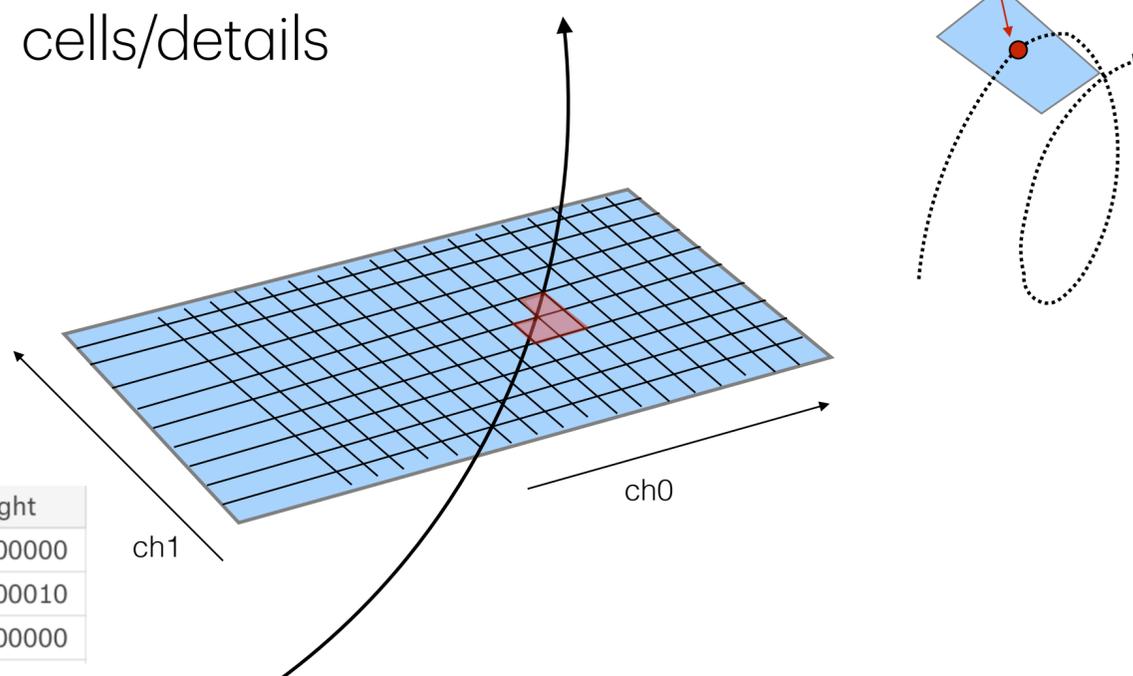


hits

	hit_id	x	y	z	volume_id	layer_id	module_id
0	1	-64.409897	-7.163700	-1502.5	7	2	1
1	2	-55.336102	0.635342	-1502.5	7	2	1
2	3	-83.830498	-1.143010	-1502.5	7	2	1
3	4	-96.109100	-8.241030	-1502.5	7	2	1
4	5	-62.673599	-9.371200	-1502.5	7	2	1
5	6	-57.068699	-8.177770	-1502.5	7	2	1
6	7	-73.872299	-2.578900	-1502.5	7	2	1
7	8	-63.853500	-10.868400	-1502.5	7	2	1
8	9	-97.254799	-10.889100	-1502.5	7	2	1
9	10	-90.292900	-3.269370	-1502.5	7	2	1
10	11	-59.182999	-0.670508	-1502.5	7	2	1

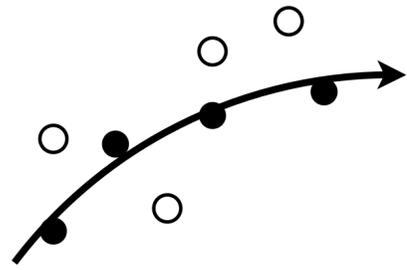
cells/details

	hit_id	ch0	ch1	value
0	1	209	617	0.013832
1	1	210	617	0.079887
2	1	209	618	0.211723
3	2	68	446	0.334087
4	3	58	954	0.034005
5	3	58	956	0.007798
6	3	60	951	0.019897



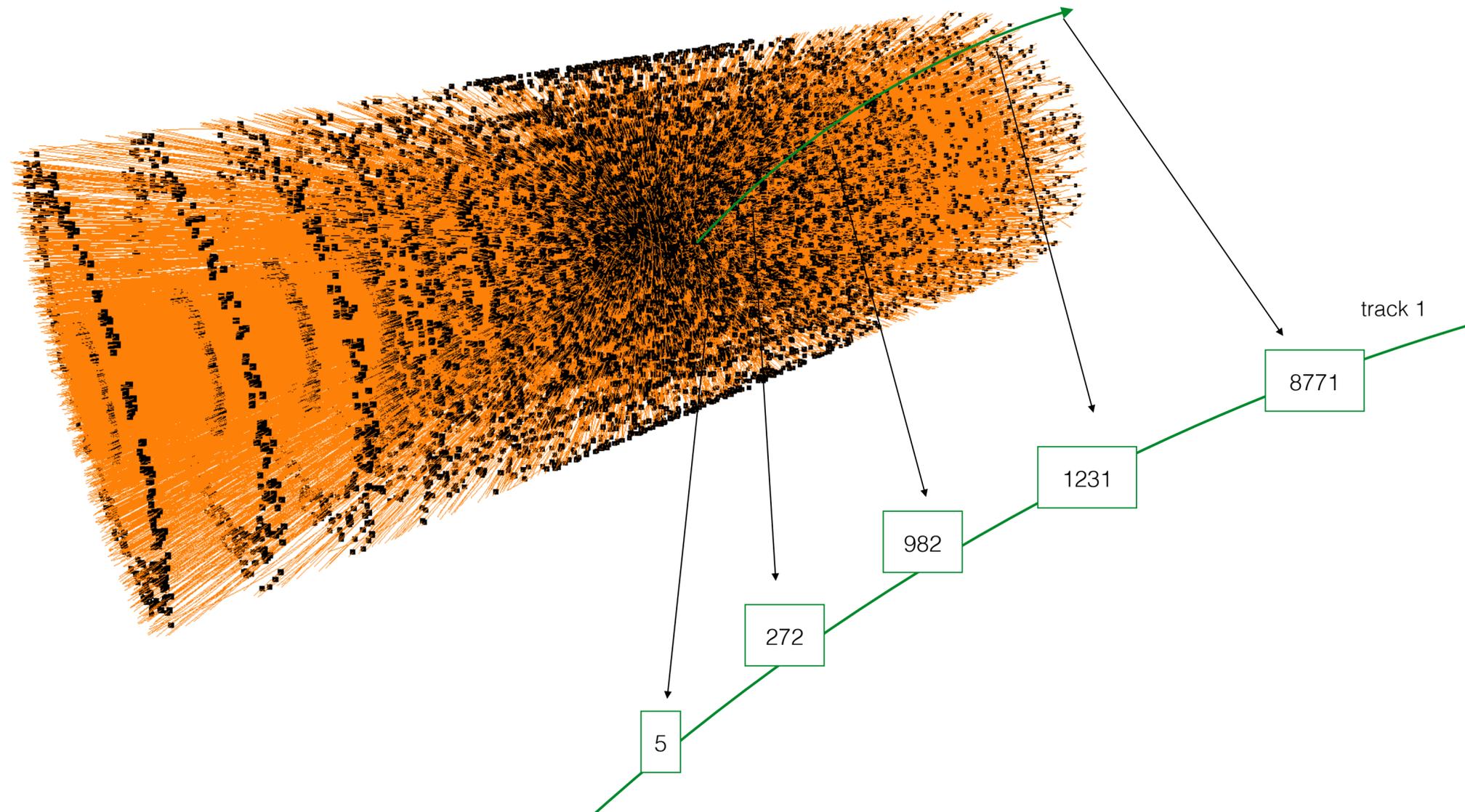
truth

	hit_id	particle_id	tx	ty	tz	tpx	tpy	tpz	weight
0	1	0	-64.411598	-7.164120	-1502.5	250710.000000	-149908.000000	-956385.000000	0.000000
1	2	22525763437723648	-55.338501	0.630805	-1502.5	-0.570605	0.028390	-15.492200	0.000010
2	3	0	-83.828003	-1.145580	-1502.5	626295.000000	-169767.000000	-760877.000000	0.000000



# The Challenge

A TrackML event - A solution

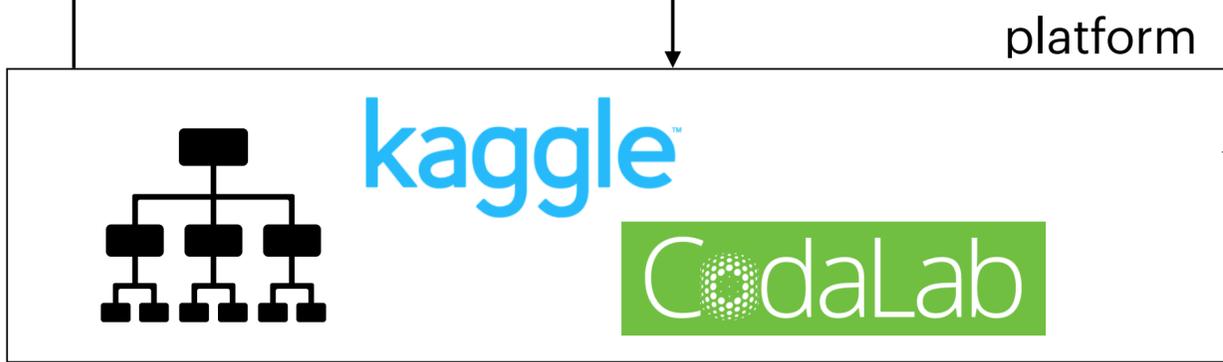
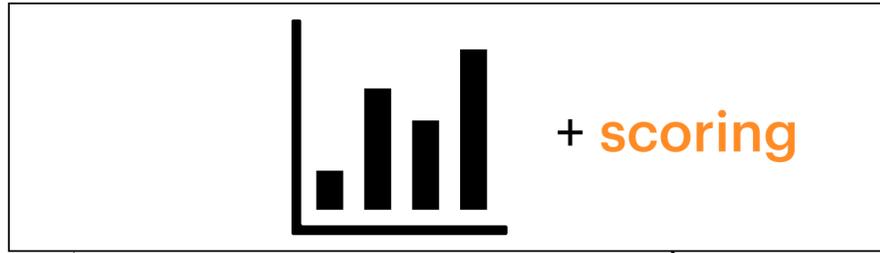


	hit_id	track_id
1	5	1
2	272	1
3	982	1
4	1231	1
5	18771	1

labelled set of hits

# Submission

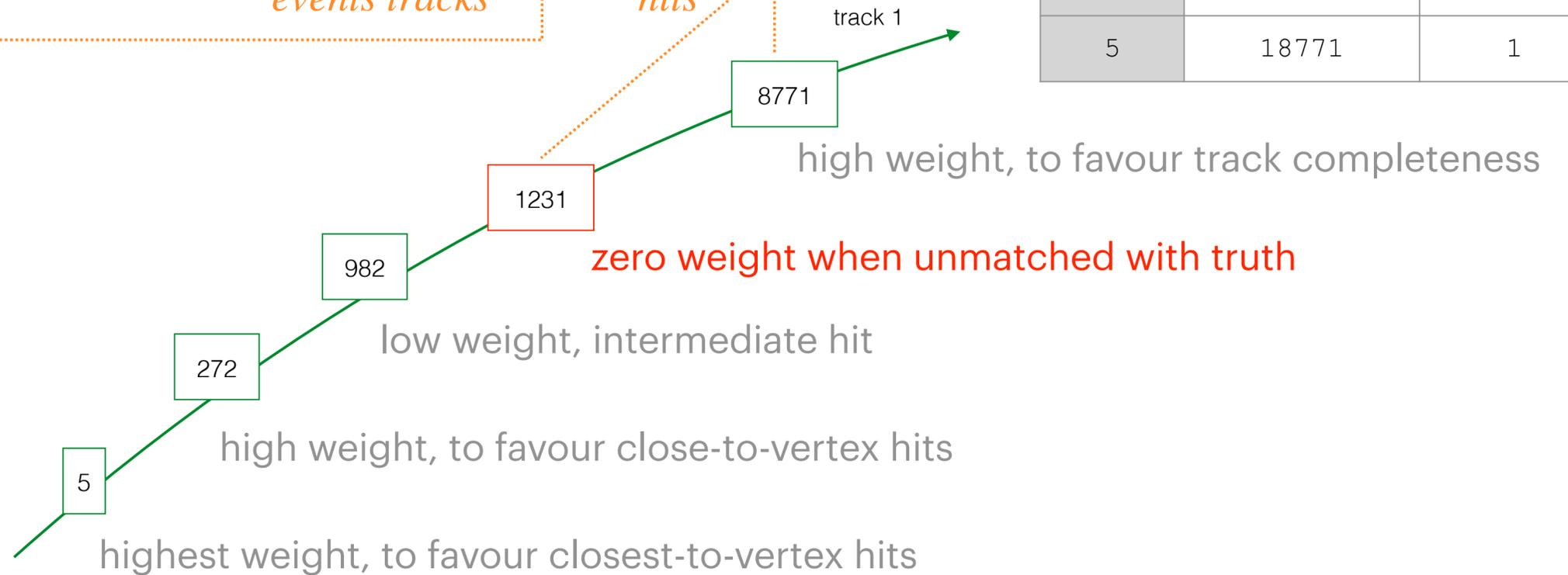
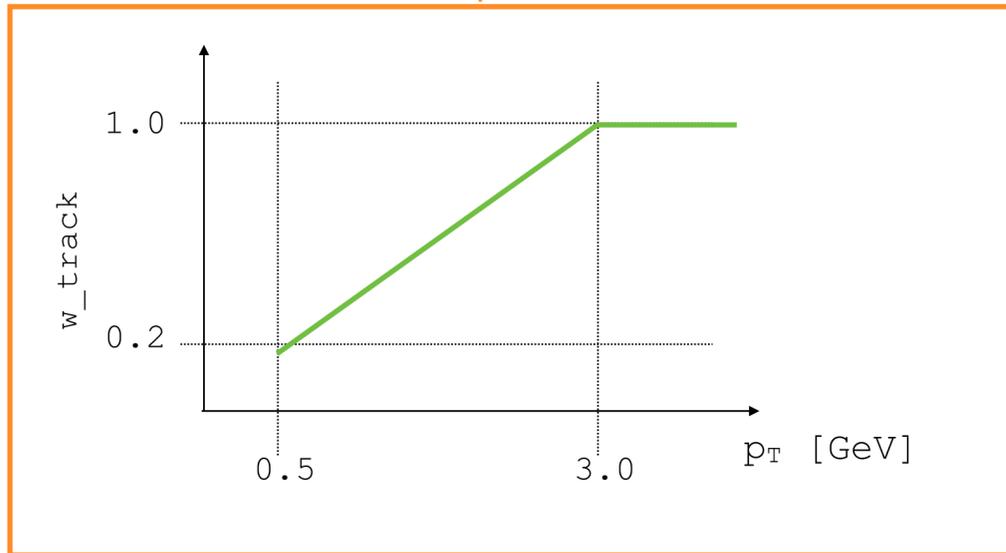
And **scoring**



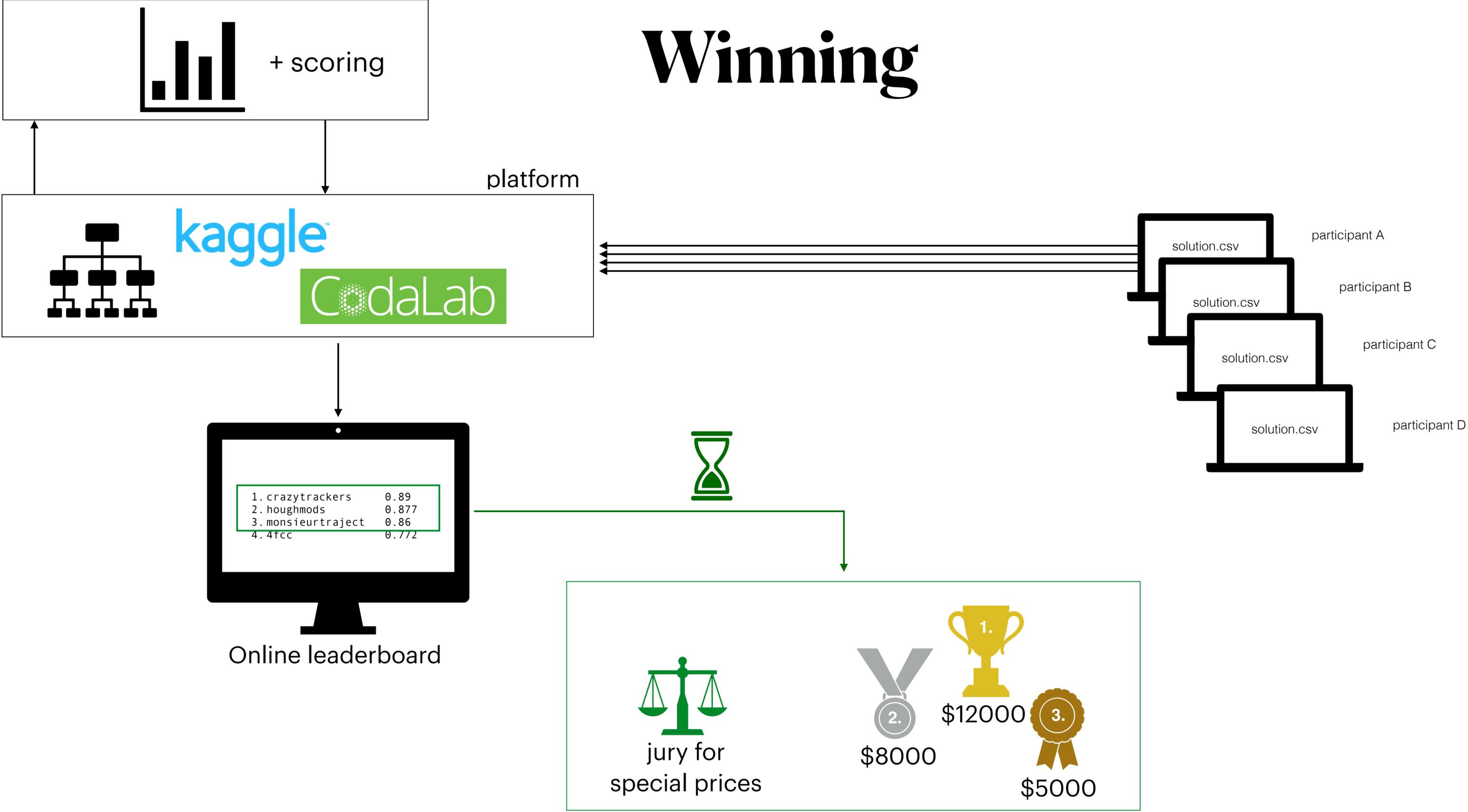
participant A

	hit_id	track_id
1	5	1
2	272	1
3	982	1
4	1231	1
5	18771	1

$$score = \sum_{events} \sum_{tracks} w_{track} \sum_{hits} w_{hit}^{matched}$$



# Winning



# Phase 1



Accuracy Challenge

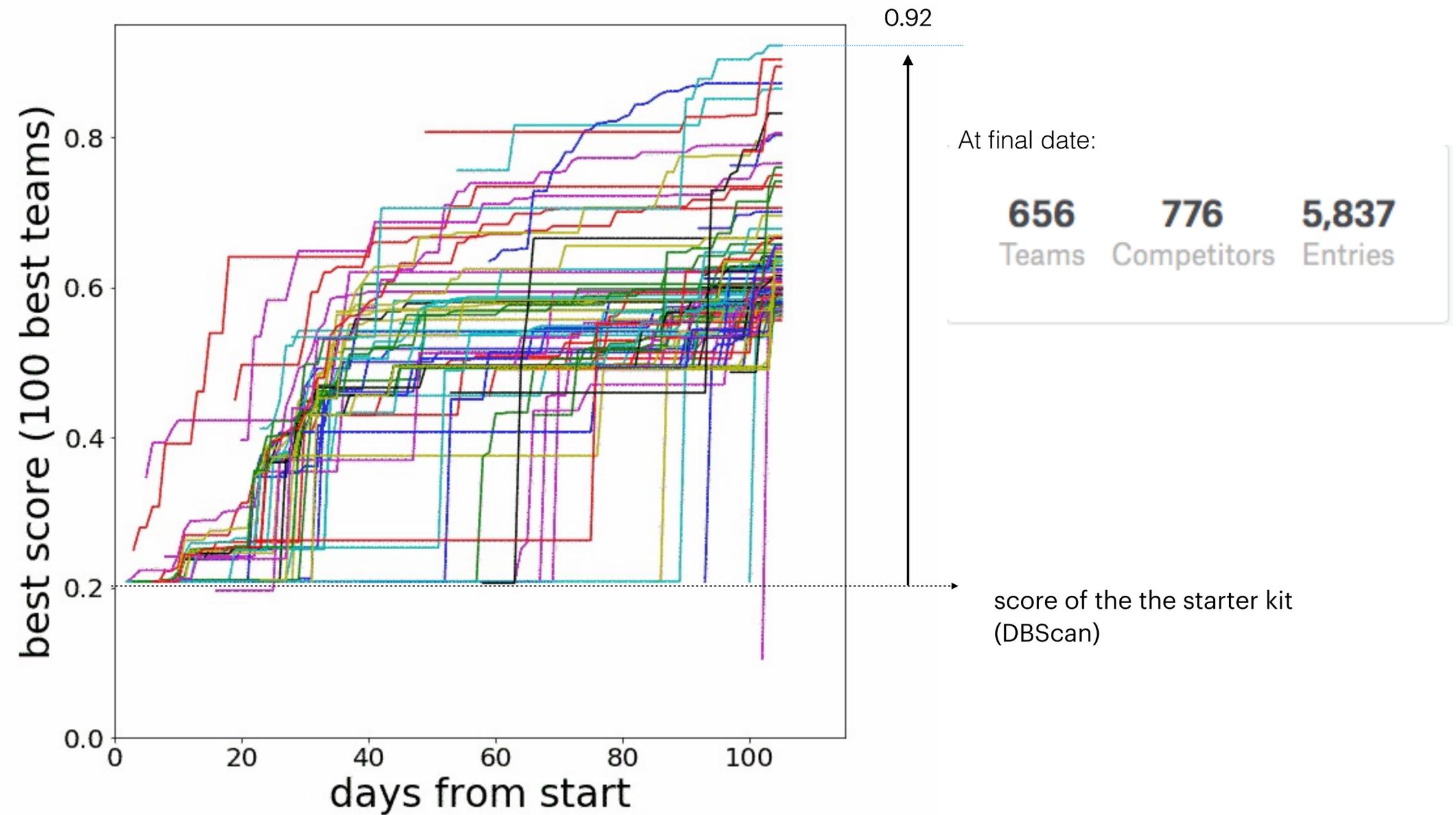
kaggle™

Apr 13, 2018



Aug 14, 2018

# Phase 1 - Score Evolution



# Phase 1 - Final Leaderboard

Public Leaderboard Private Leaderboard

The private leaderboard is calculated with approximately 71% of the test data.  
This competition has completed. This leaderboard reflects the final standings. [Refresh](#)

■ In the money
 ■ Gold
 ■ Silver
 ■ Bronze

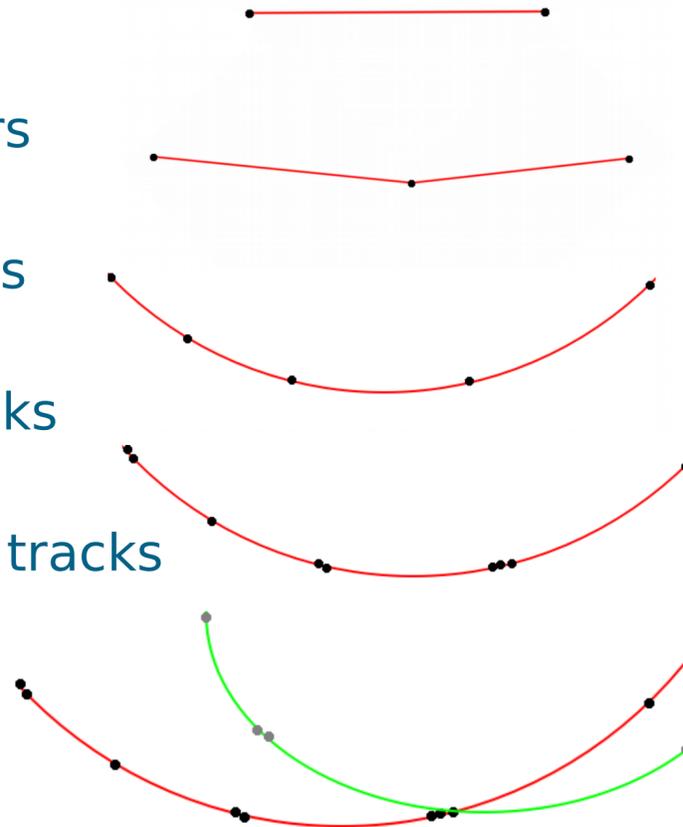
#	Δpub	Team Name	Kernel	Team Members	Score	Entries	Last
1	—	Top Quarks			0.92182	10	2mo
2	—	outrunner			0.90302	9	2mo
3	—	Sergey Gorbunov			0.89353	6	2mo
4	—	demelian			0.87079	35	2mo
5	—	Edwin Steiner			0.86395	5	2mo
6	—	Komaki			0.83127	22	2mo
7	—	Yuval & Trian			0.80414	56	2mo
8	—	bestfitting			0.80341	6	2mo

# Phase 1 Top Quarks



## Main steps

- Select promising pairs
  - 7 million / 0.99
- Extend pairs to triples
  - 12 million / 0.97
- Extend triples to tracks
  - 12 million / 0.95
- Add duplicate hits to tracks
  - 12 million / 0.96
- Assign hits to tracks
  - 90% of hits / 0.92



## Findings

- No magic formula
- We won because we were fast to try out and implement many ideas and got the details right
  - I once earned 0.03 (0.85→0.88) from fixing a tuning parameter
- In other words: combination of many factors

- Logistic regression for track candidate pruning

- Pure C++, some scikit-learn for training



Author	Johan S. Wind
<Wall time>/evt	7m 17s
Peak memory	2.78 Gb

# Phase 1 Outrunner



- **ML approach using python & keras**

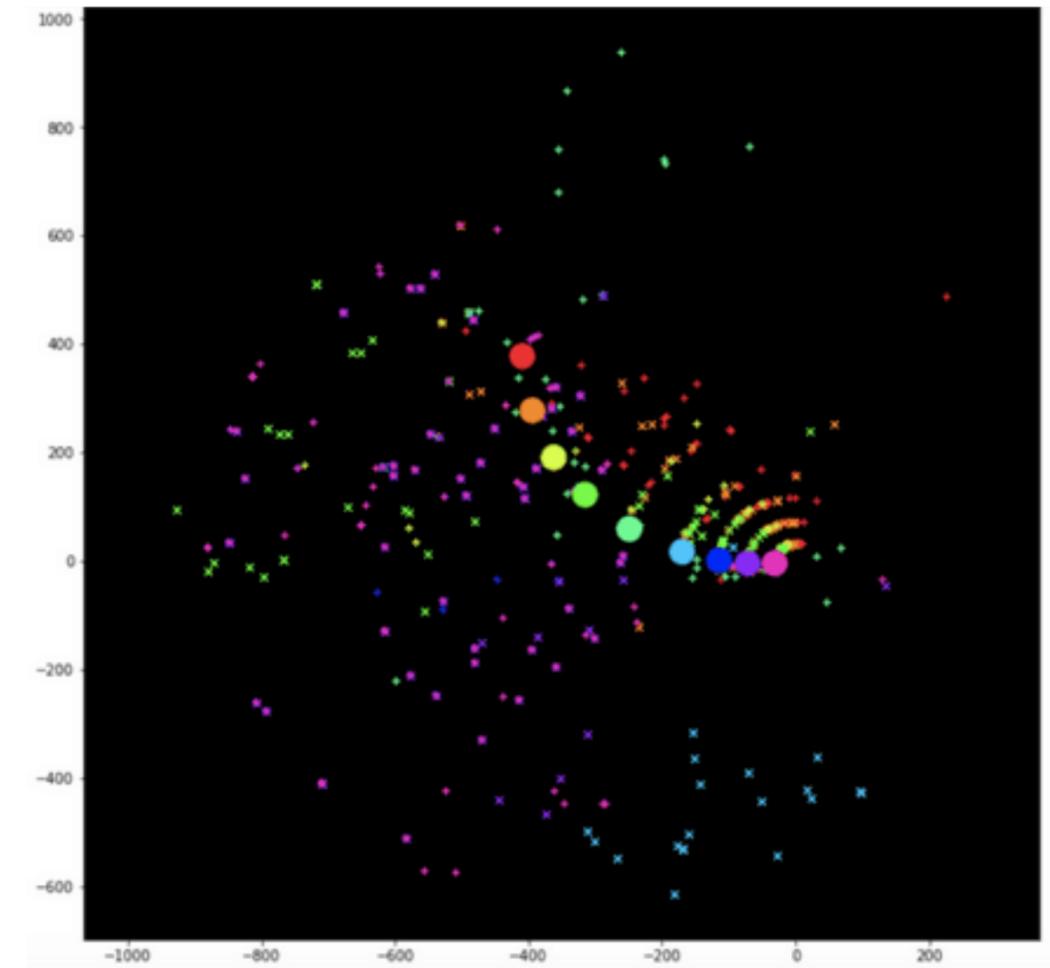
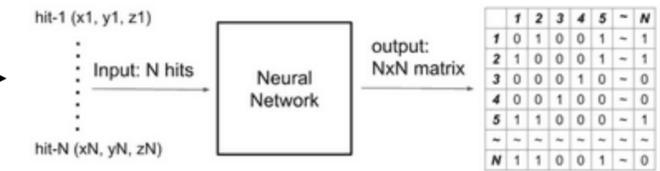
- event with N hits
- predict N x N relationships between hits, connect pairs when their probability is 1

- **Training**

- 5 hidden layers with 4k - 2k - 2k - 1k nodes
- 27 input features

- **Prediction & Reconstruction**

- predict pairwise relationship probability
- finding highest probability pair & add pairwise



 PY	Author	Pei-Lien Chou
	<Wall time>/evt	~1 day

# Phase 1 Outrunner



- **ML approach using python & keras**

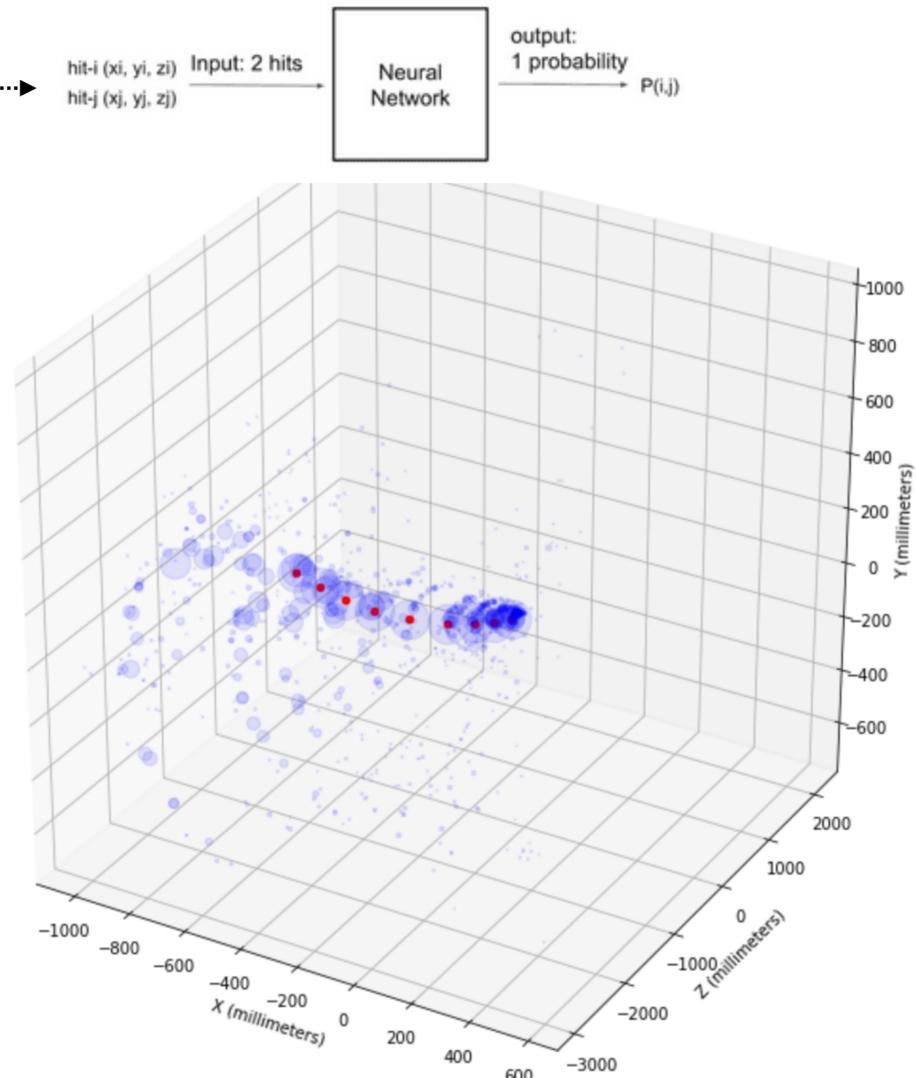
- event with N hits
- predict N x N relationships between hits, connect pairs when their probability is 1

- **Training**

- 5 hidden layers with 4k - 2k - 2k - 1k nodes
- 27 input features

- **Prediction & Reconstruction**

- predict pairwise relationship probability
- finding highest probability pair & add pairwise

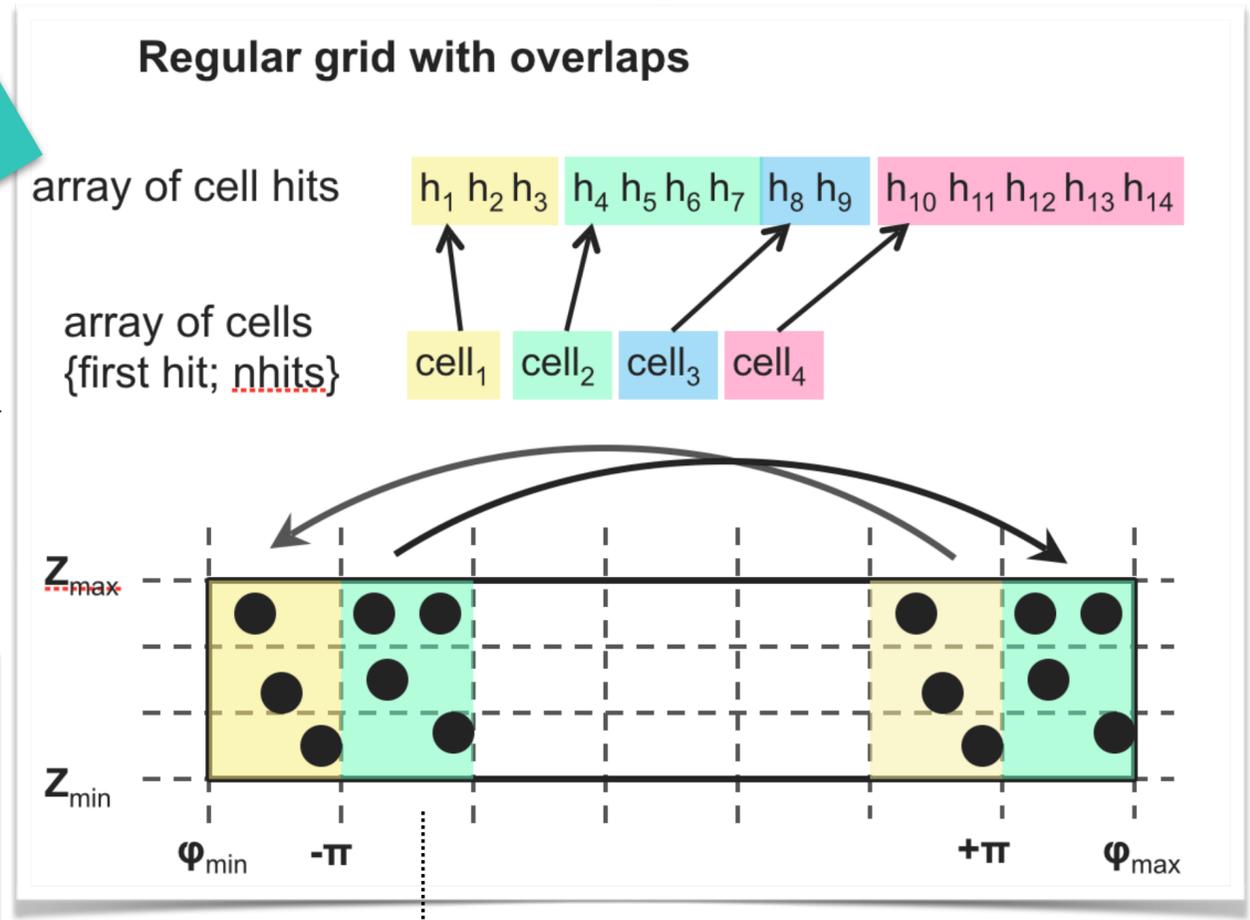
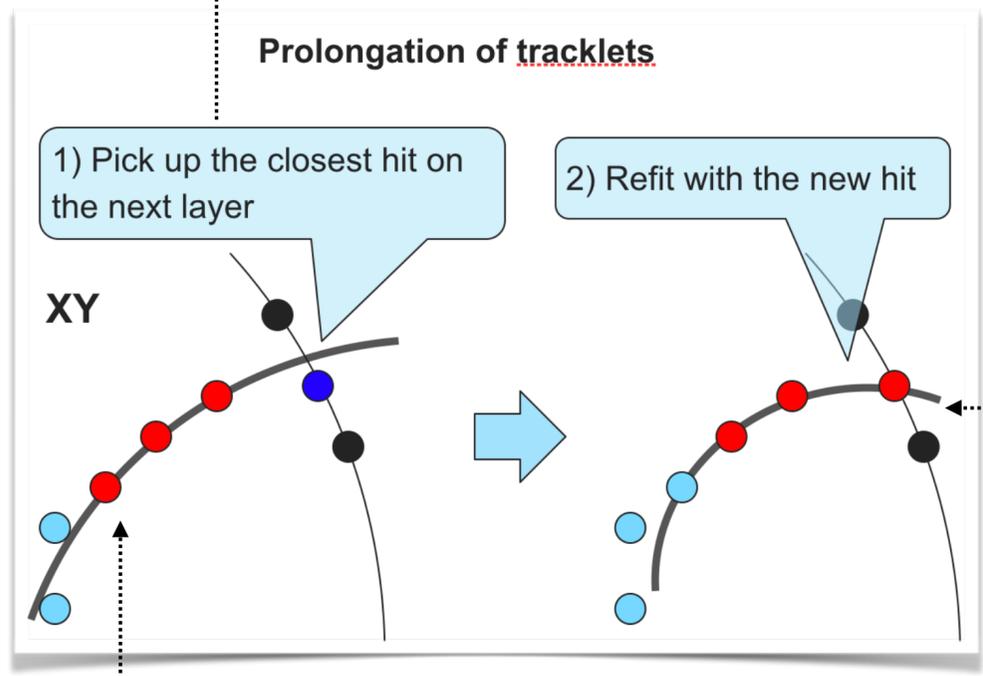
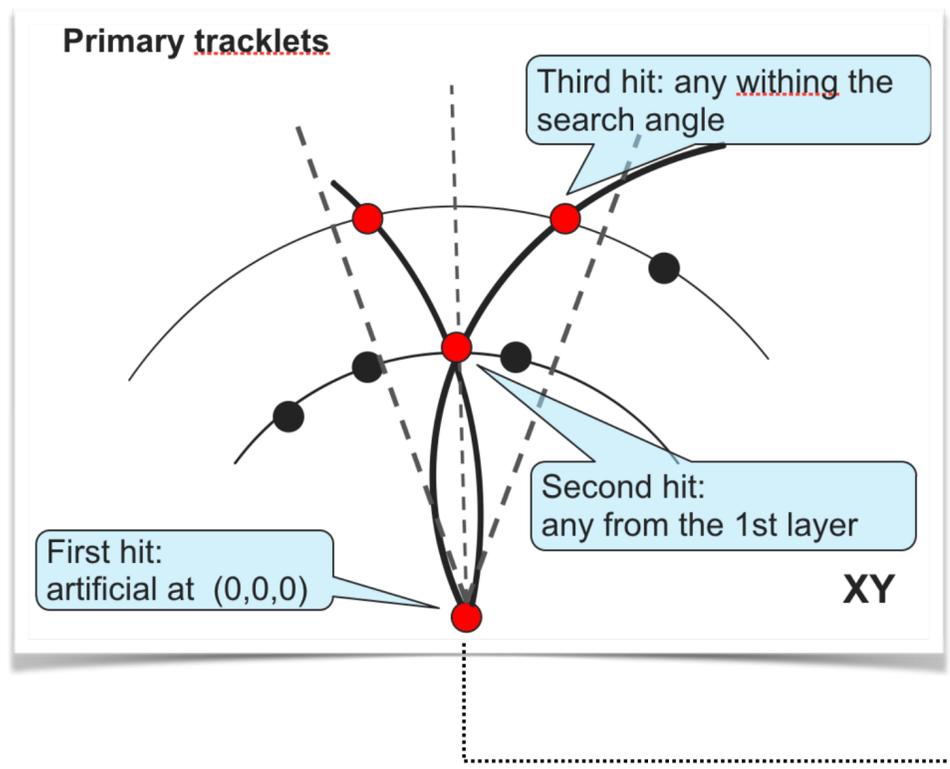


	Author	Pei-Lien Chou
	<Wall time>/evt	~1 day

# Phase 1 Sergey Gorbunov



- Combinatorial approach based on track following
- No search branches
- Simple track model (Helix)



	Author	[ <a href="#">Sergey Gorbunov</a> ]
	<Wall time>/evt	~1.2 mins

# Phase 1 Jury Prizes

## Innovation prize

Yuval Reina & Trian Xylouris  
Marginalized Hough transform with machine learning classifier

## Clustering prize

Jean-Francois Puget (kaggle grandmaster)  
DBScan clustering with iterative Hough transform

## Deep Learning prize

Nicole & Liam Finnie  
DBScan seeding and LSTM track Building

## Organizer's pick

Diogo R. Ferreira  
Innovative pattern matching

#	△pub	Team Name	Kernel	Team Members	Score	Entries
1	—	Top Quarks			0.92182	10
2	—	outrunner	<b>In the money</b>		0.90302	9
3	—	Sergey Gorbunov			0.89353	6
4	—	demelian			0.87079	35
5	—	Edwin Steiner			0.86395	5
6	—	Komaki			0.83127	22
7	—	Yuval & Trian	<b>Jury pick</b>		0.80414	56
8	—	bestfitting			0.80341	6
9	—	DBSCAN forever	<b>Jury pick</b>		0.80114	23
10	—	Zidmie & KhaVo			0.76320	26
11	—	Andrea Lonza			0.75845	15
12	—	Finnies	<b>Jury pick</b>		0.74827	56
13	—	Rei Matsuzaki			0.74035	12
14	—	Mickey			0.73217	10
15	—	Vicens Gaitan			0.70429	19
16	—	Robert			0.69955	3

100	▲ 2	Diogo	<b>Jury pick</b>		0.55480	22	8mo
-----	-----	-------	------------------	--	---------	----	-----

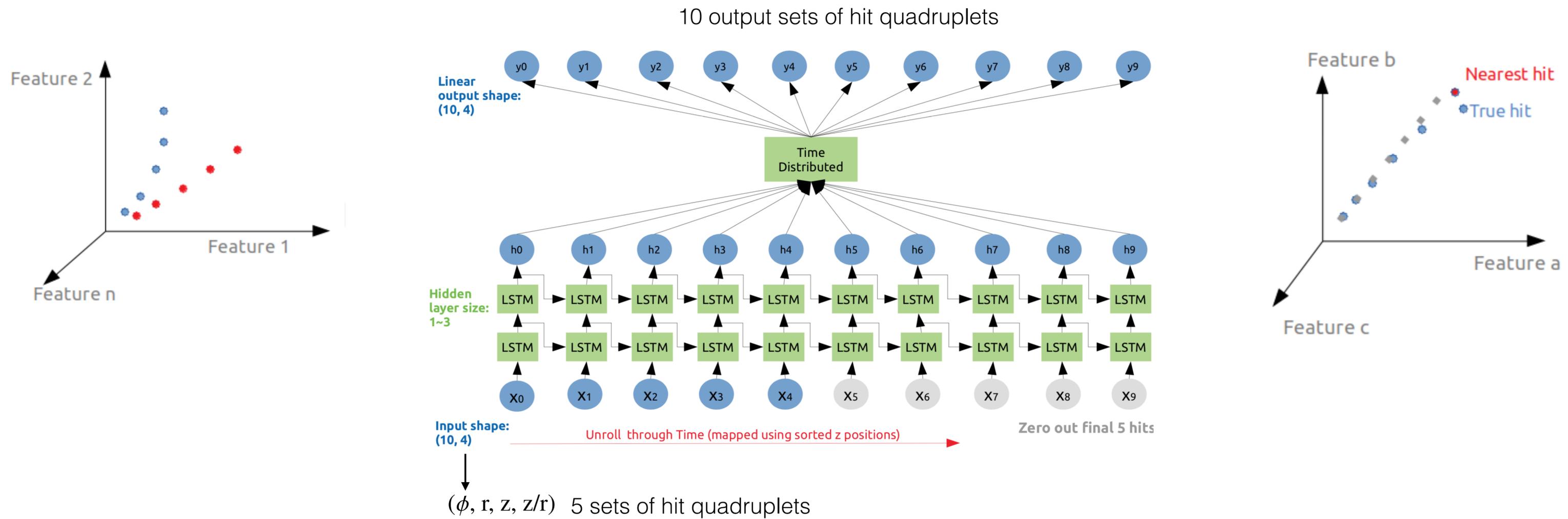
# Phase 1 DL Prize



Seeding  
DBScan

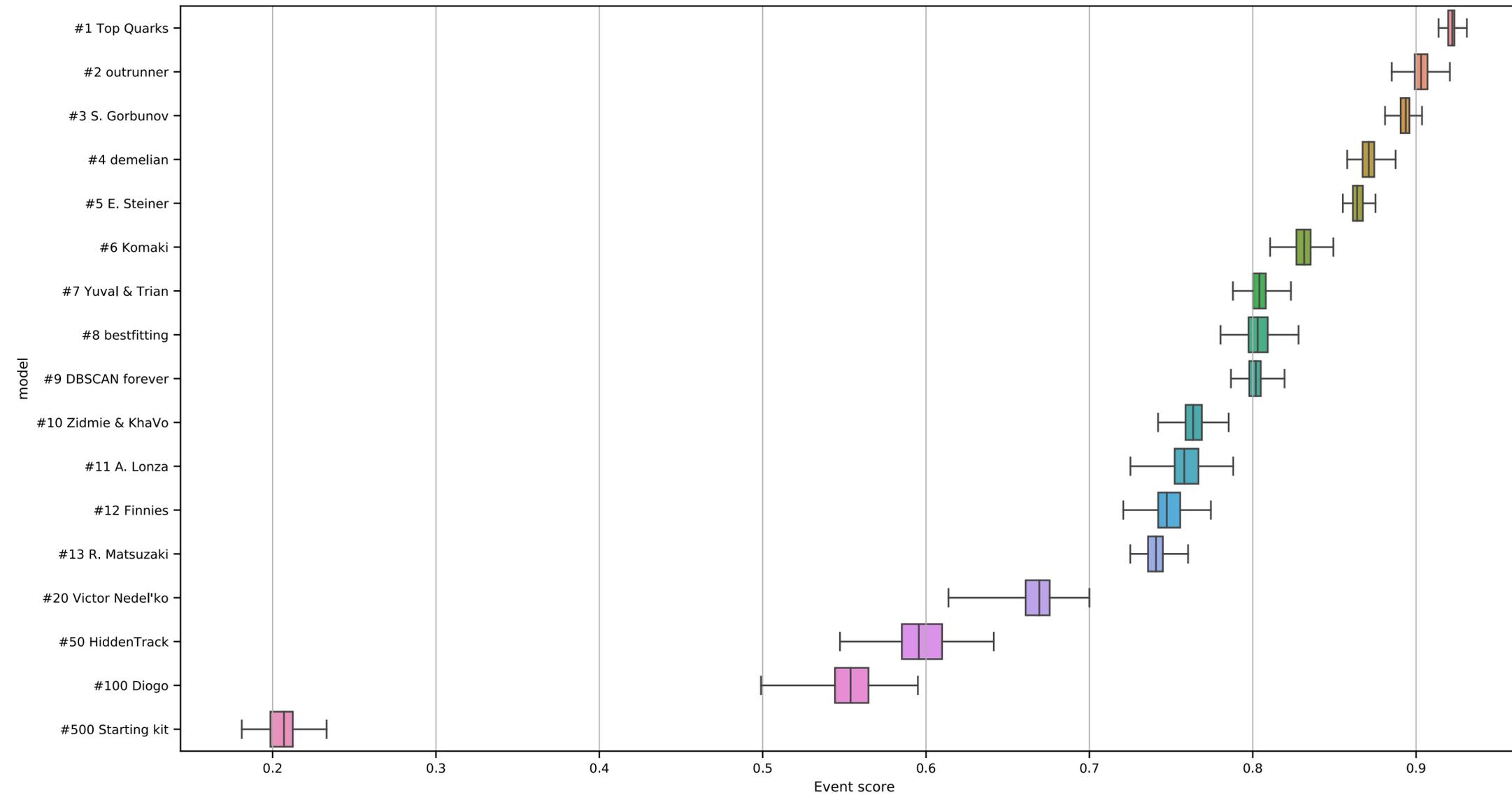
Inference & Assembling  
Recurrent NN

Fitting  
KNN-tree



# Phase 1 Aftermath - Score Stability

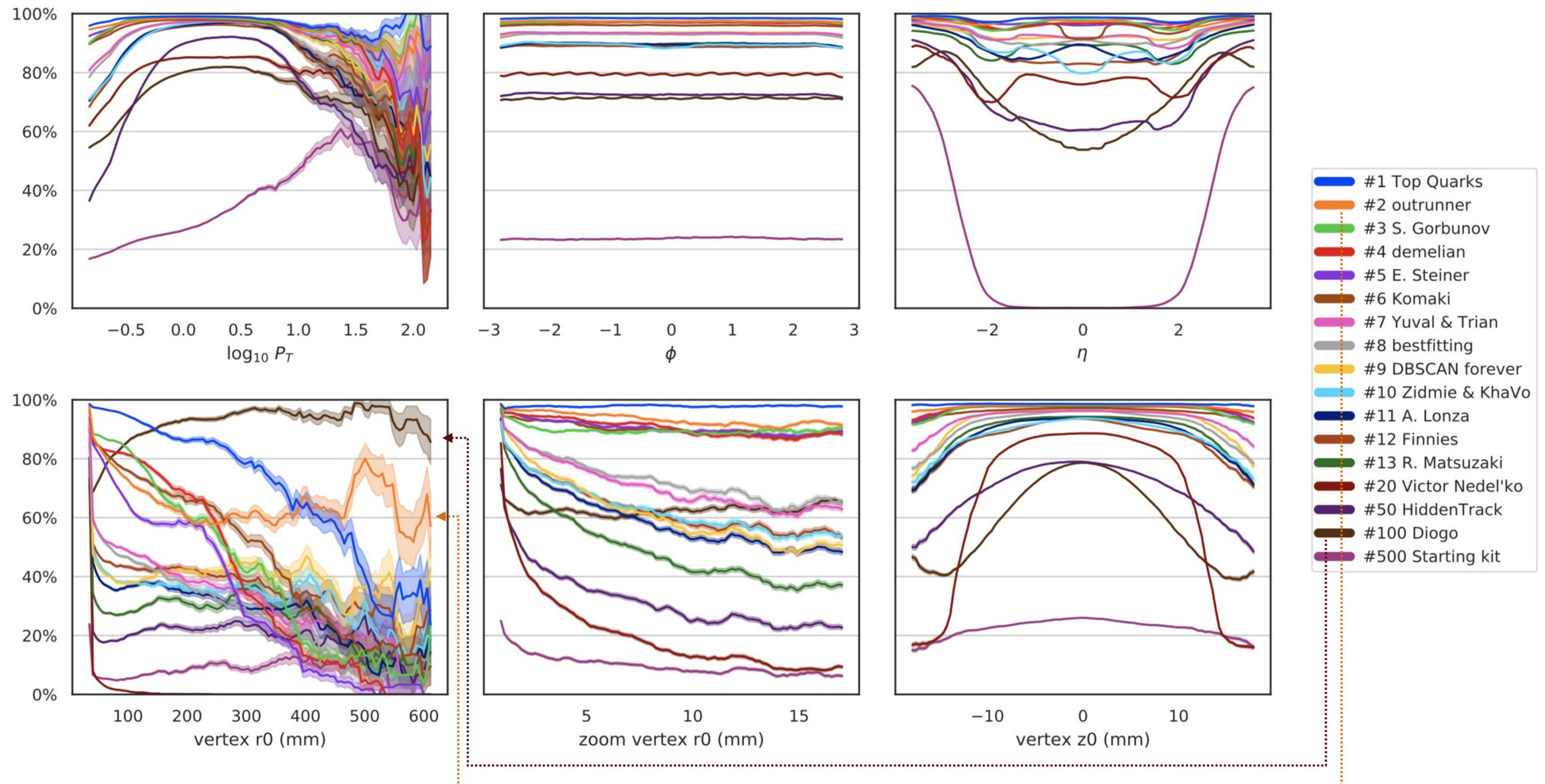
Score is excellent **measure of performance.**



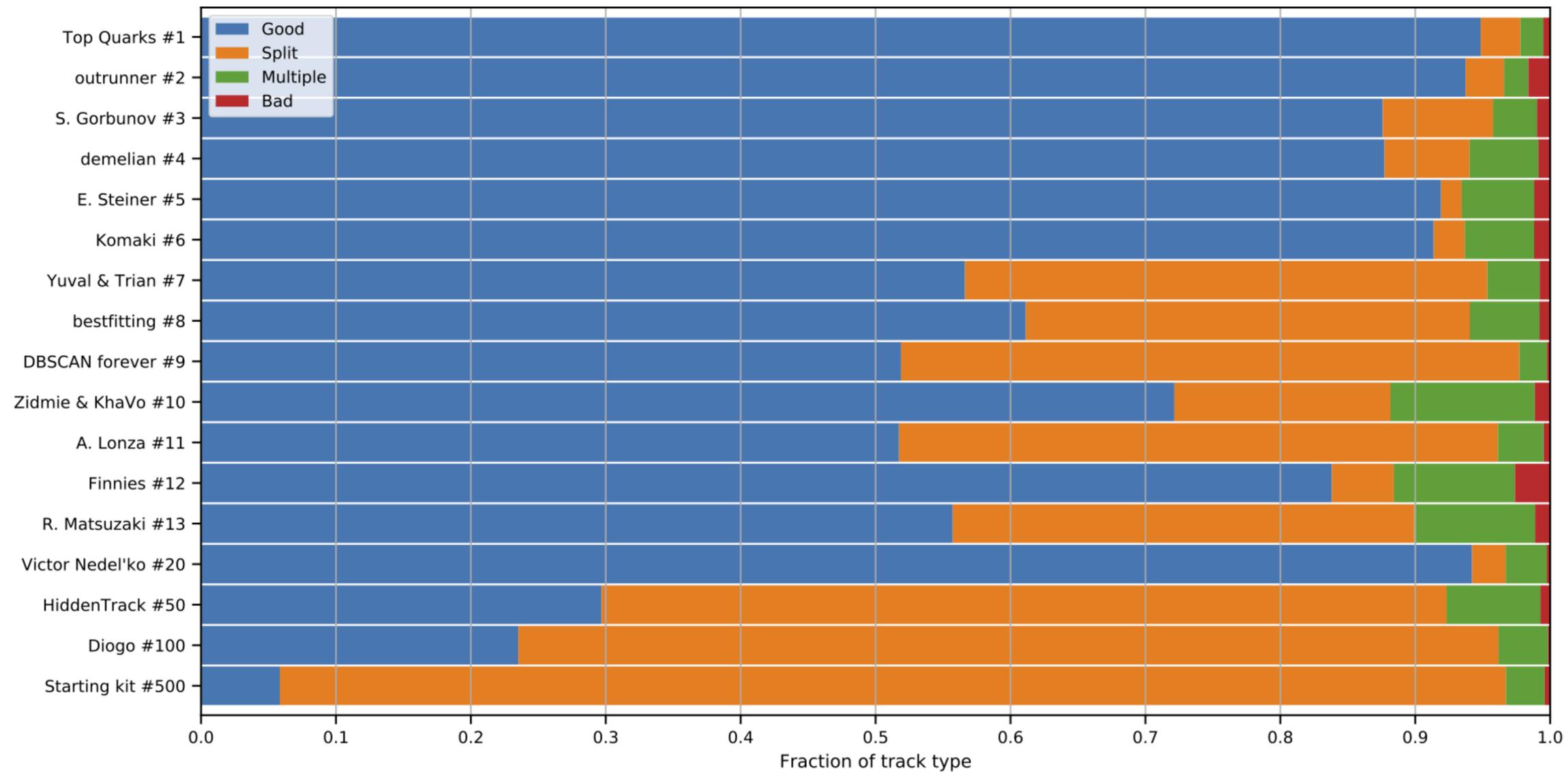
Score quartiles/extrema of submitted solutions.

# Phase 1 Aftermath - Tracking Efficiency

High score means High tracking efficiency



# Phase 1 Aftermath - Track classes



**Good:** track and particle purities above 50% (goes into the score)

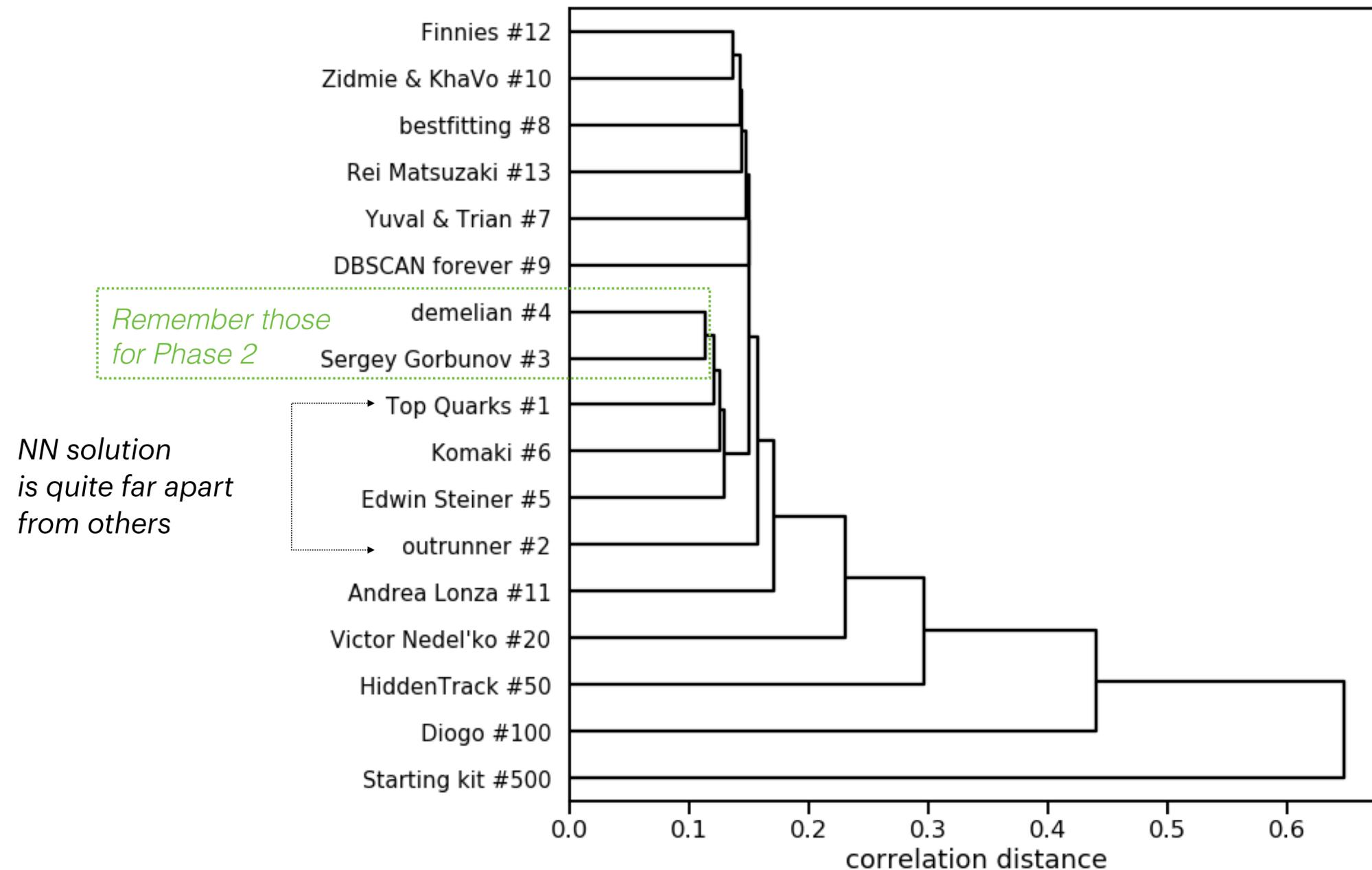
**Split:** particle purity below 50%, track purity above 50%

**Multiple:** particle purity above 50%, but track purity below 50%

**Bad:** both below 50%

# Phase 1 Aftermath - Score correlations

High score means High tracking efficiency



# Phase 2



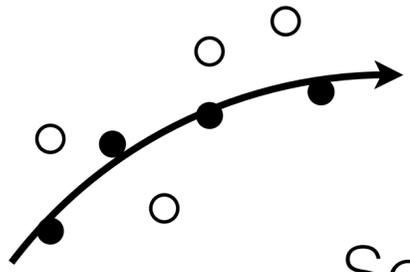
Accuracy & Speed Challenge



Sep 07, 2018

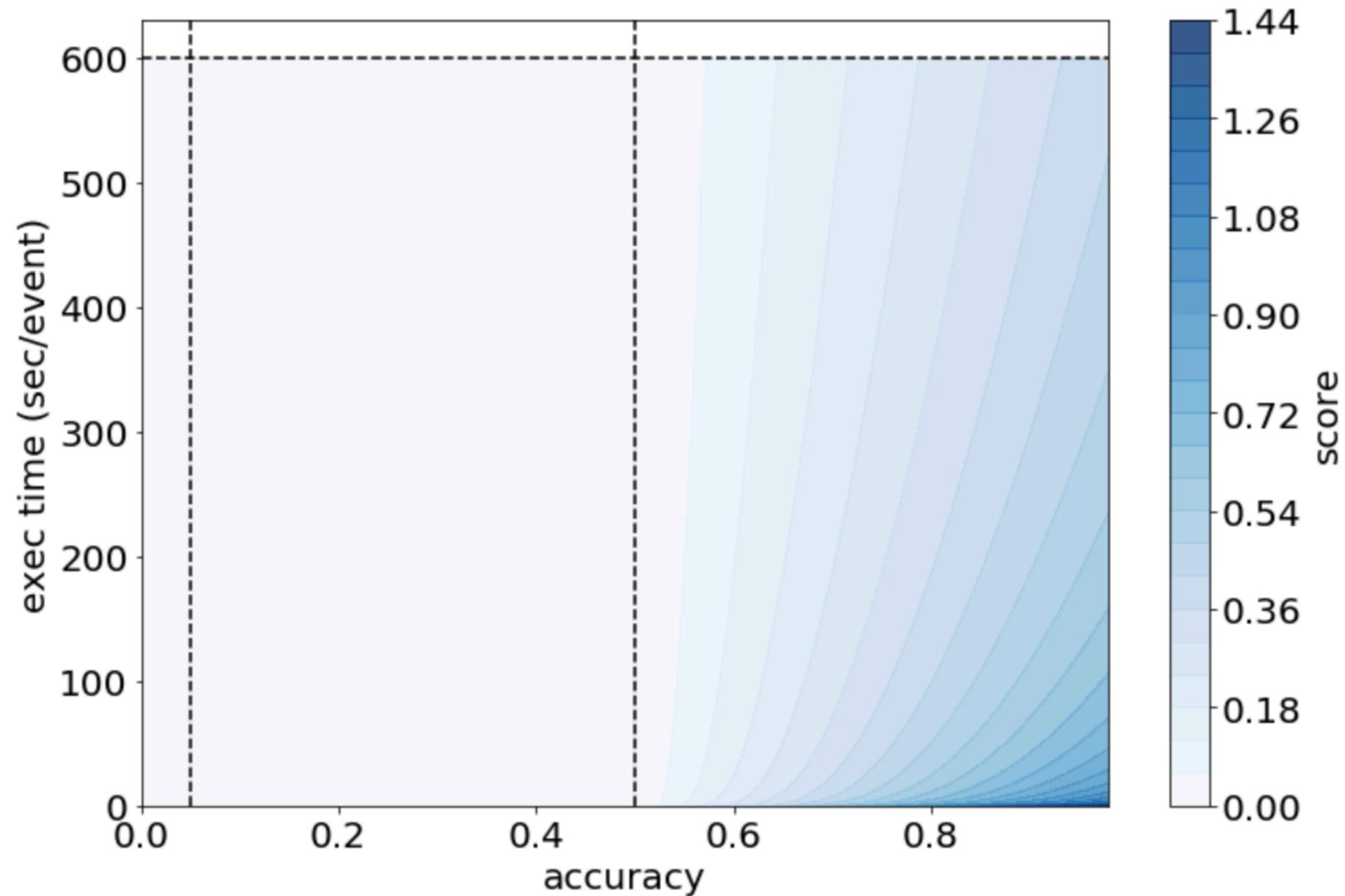


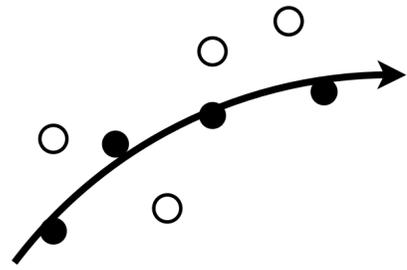
Nov 12, 2018



# Phase 2 - Modifications

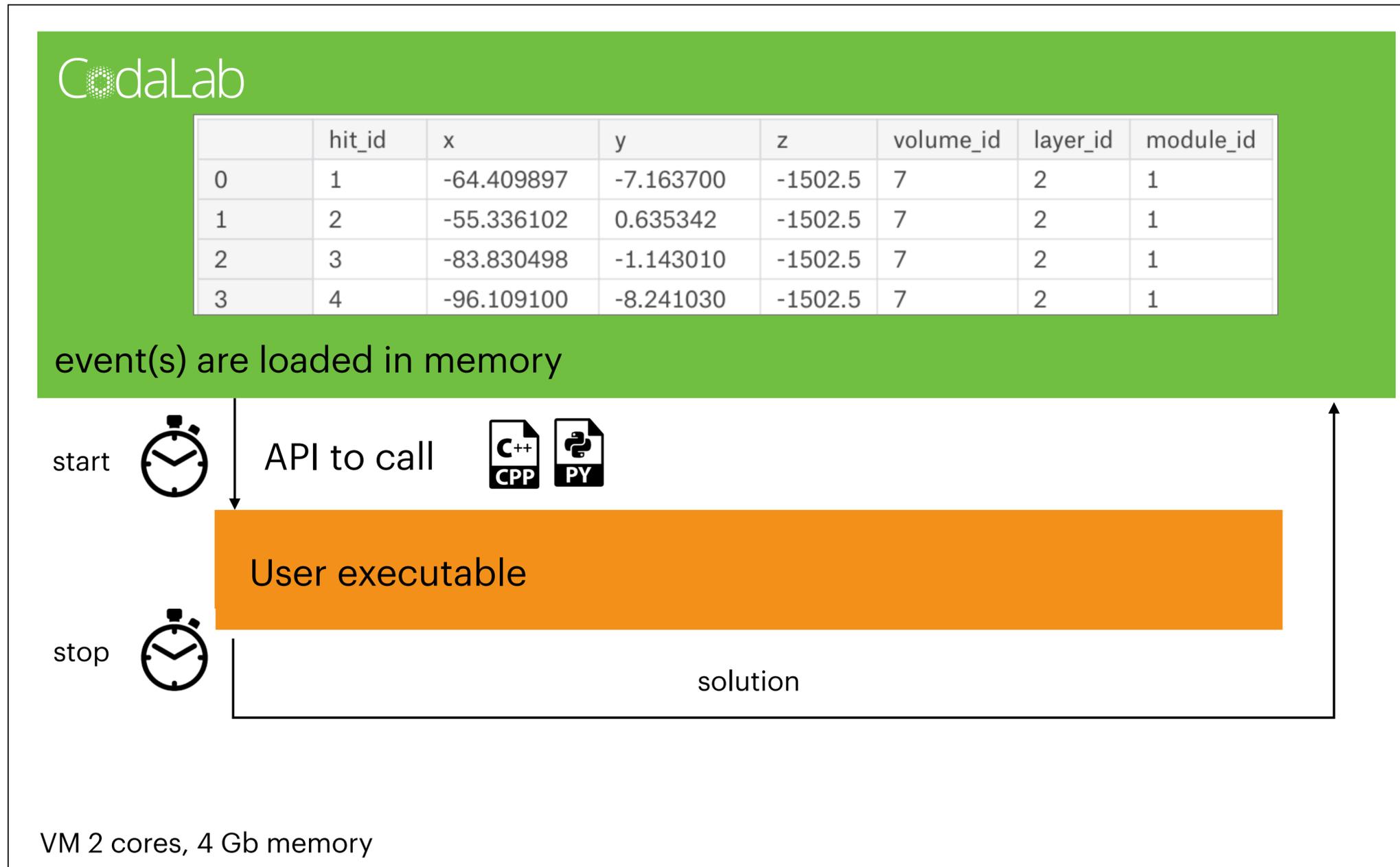
Score built from two components + restricted to primary particles only



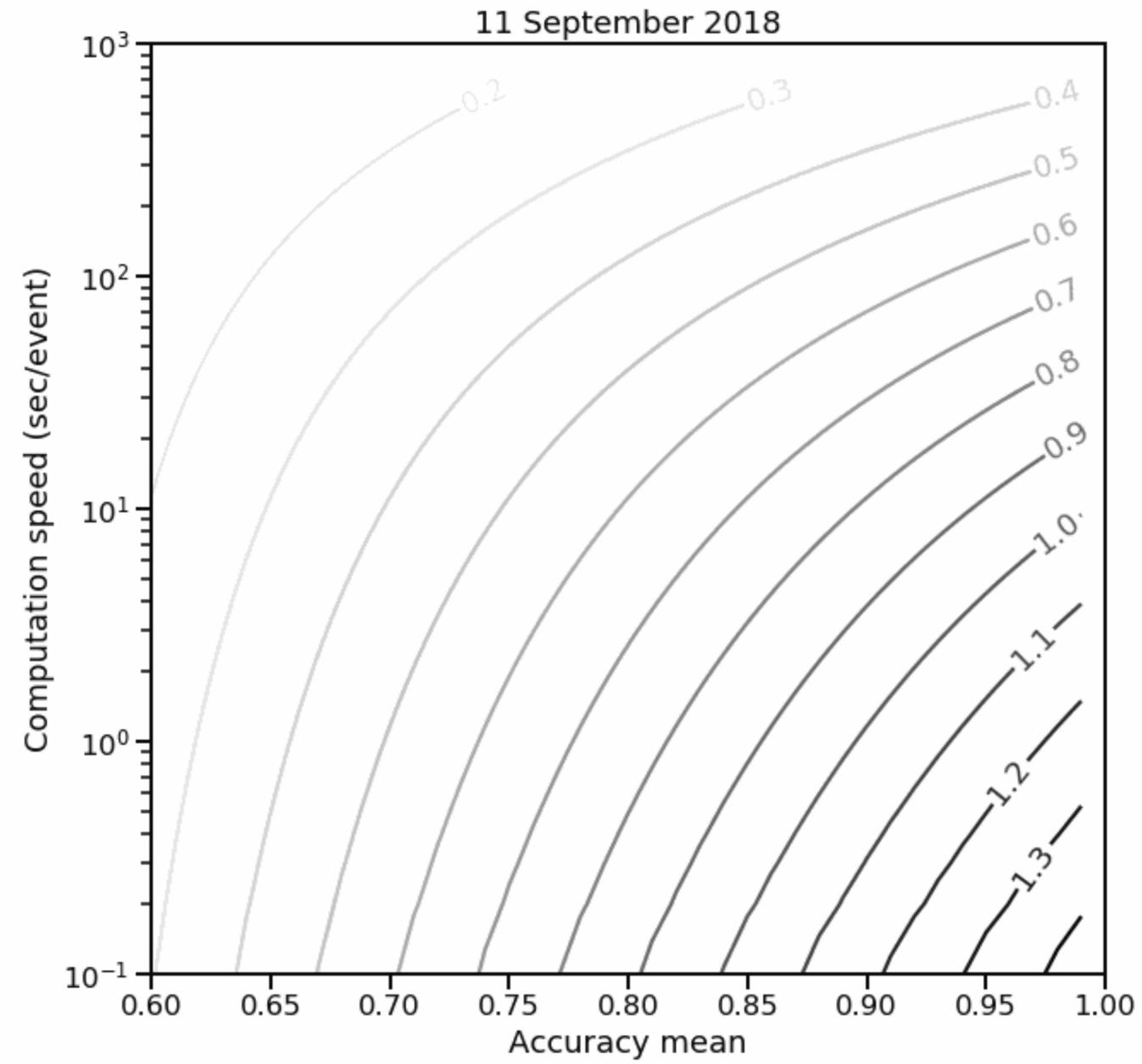


# Phase 2 - Modifications

Runtime control



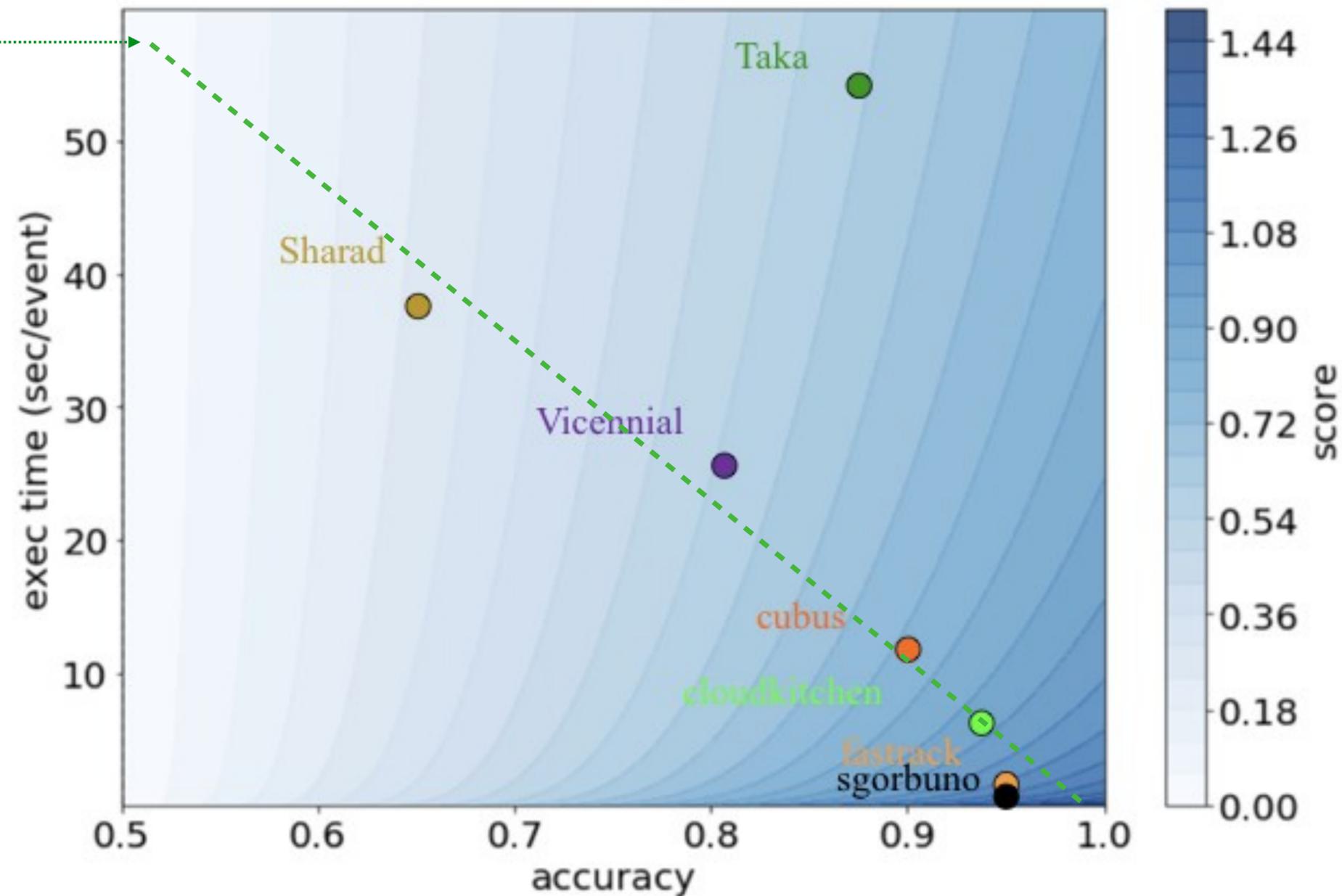
# Phase 2 - Score Evolution



# Phase 2 - Final Score Map

Correlation accuracy / speed

Fastest solution are  
in general also the most  
accurate one!



# Phase 2 - Final Leaderboard

Only 7 teams made the cut!

RESULTS										
#	User	Entries	Date of Last Entry	score ▲	accuracy_mean ▲	accuracy_std ▲	computation time (sec) ▲	computation speed (sec/event) ▲	Duration ▲	
1	<b>sgorbuno</b> 	9	03/12/19	1.1727 (1)	0.944 (2)	0.00 (14)	28.06 (1)	0.56 (1)	64.00 (1)	
2	<b>fastrack</b> 	53	03/12/19	1.1145 (2)	0.944 (1)	0.00 (15)	55.51 (16)	1.11 (16)	91.00 (6)	
3	<b>cloudkitchen</b> 	73	03/12/19	0.9007 (3)	0.928 (3)	0.00 (13)	364.00 (18)	7.28 (18)	407.00 (8)	
4	cubus	8	09/13/18	0.7719 (4)	0.895 (4)	0.01 (9)	675.35 (19)	13.51 (19)	724.00 (9)	
5	Taka	11	01/13/19	0.5930 (5)	0.875 (5)	0.01 (12)	2668.50 (23)	53.37 (23)	2758.00 (13)	
6	Vicennial	27	02/24/19	0.5634 (6)	0.815 (6)	0.01 (10)	1270.73 (20)	25.41 (20)	1339.00 (10)	
7	Sharad	57	03/10/19	0.2918 (7)	0.674 (7)	0.02 (4)	1902.20 (22)	38.04 (22)	1986.00 (12)	
8	WeizmannAI	5	03/12/19	0.0000 (8)	0.133 (11)	0.01 (11)	88.08 (17)	1.76 (17)	124.00 (7)	
9	harshakoundinya	2	03/12/19	0.0000 (8)	0.085 (13)	0.01 (6)	49.22 (8)	0.98 (8)	86.00 (3)	
10	iWit	6	03/10/19	0.0000 (8)	0.082 (15)	0.01 (8)	48.23 (3)	0.96 (3)	85.00 (2)	

# Phase 2 Sergey Gorbunov



- “Mikado” algorithm

- Based on 3rd solution of Phase 1

- Modifications

- Runs iteratively in 80 passes with hit removal
- Search branches are enabled
- Parameters for each pass optimised (semi-automated):

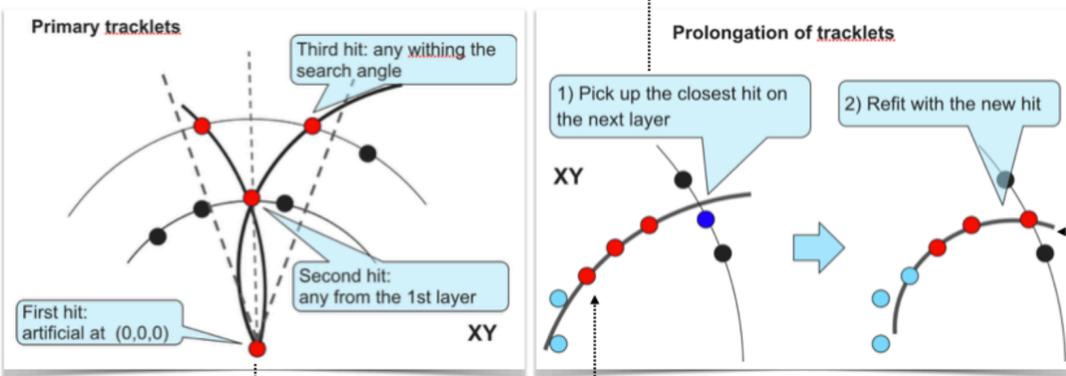
$O(10k)$  parameters optimized

- High Accuracy

- Reaching 0.944 (primary particles only)

## Phase 1 Sergey Gorbunov

- Combinatorial approach based on track following
- No search branches
- Simple track model (Helix)



**Primary tracklets**

First hit: artificial at (0,0,0)

Second hit: any from the 1st layer

Third hit: any within the search angle

**Prolongation of tracklets**

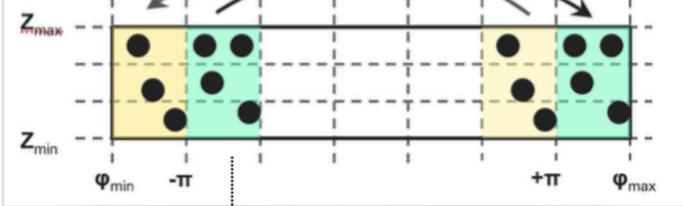
1) Pick up the closest hit on the next layer

2) Refit with the new hit

**Regular grid with overlaps**

array of cell hits:  $h_1 h_2 h_3 h_4 h_5 h_6 h_7 h_8 h_9 h_{10} h_{11} h_{12} h_{13} h_{14}$

array of cells {first hit; nhits}:  $cell_1 cell_2 cell_3 cell_4$



	Author	[ <a href="#">Sergey Gorbunov</a> ]
	<Wall time>/evt	~1.2 mins

	Author	[ <a href="#">Sergey Gorbunov</a> ]
	<Wall time>/evt	0.56s
	Peak memory	0.1 Gb

# Phase 2 FASTTrack



- **Iterative algorithm**

- Based on 4th solution of Phase 1

- **Algorithm outline**

- using **measurement shapes** to predict track inclination
- segment based track following
  - Using connection graph pre-built from Detector.csv
  - Run Cellular Automaton (CA), OpenMP parallelized
  - Candidate building: graph traversal with applied KF
  - Combinatorial kalman filtering

- **High Accuracy**

- Reaching **0.944** (primary particles only)

Iteratively done 3 times  
with hit removal

	Author	[ <a href="#">Sergey Gorbunov</a> ]
	<Wall time>/evt	1.1 s
	Peak memory	0.6 Gb

+ OpenMP

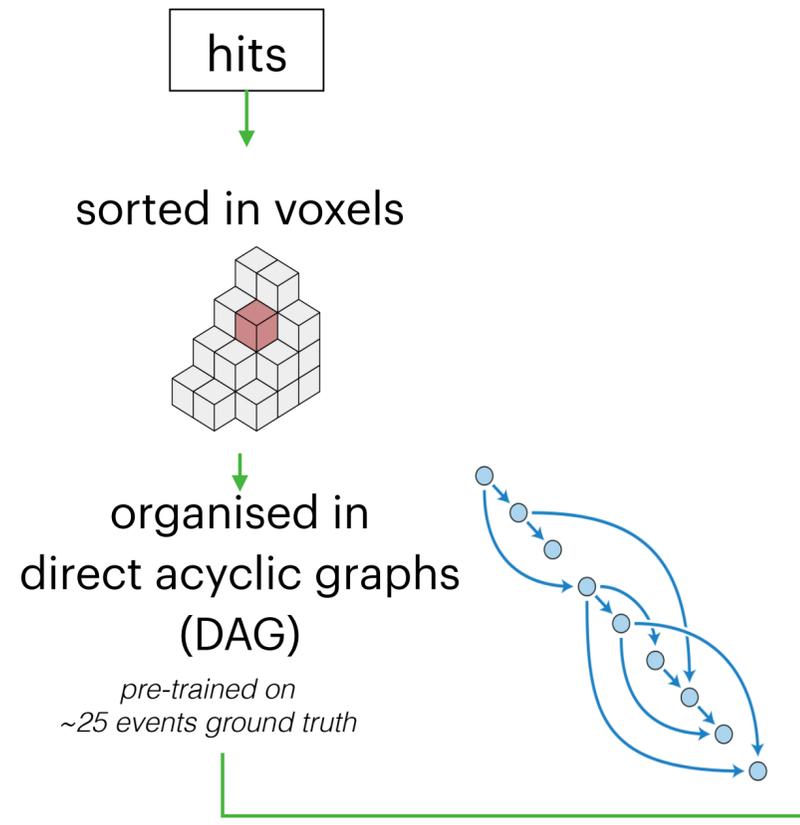
# Phase 2 Cloud kitchen



- Partly reused concept

- Based on winning solution of Phase 1

- Achieved accuracy: **0.93**



## Phase 1 Top Quarks

### Main steps

- Select promising pairs
  - 7 million / 0.99
- Extend pairs to triples
  - 12 million / 0.97
- Extend triples to tracks
  - 12 million / 0.95
- Add duplicate hits to tracks
  - 12 million / 0.96
- Assign hits to tracks
  - 90% of hits / 0.92

### Findings

- No magic formula
- Triplet finder

Logistic regression for track candidate pruning

Pure C++, some scikit-learn for training

Doublet finder

NN1

NN2

Triplet finder

NN3

	Author	Johan S. Wind
	<Wall time>/evt	7m 17s
	Peak memory	2.78 Gb

$\pm z$  graph set  
 $\eta - \phi$  graph set

Threaded

	Author	[ <a href="#">Marcel Kunze</a> ]
	<Wall time>/evt	~7 s
	Peak memory	0.79 Gb

# Key Lessons

## Phase 1

- **Huge interest in/outside community**
  - Reputation of kaggle platform helped
- **Worth investing in a good score**
  - Solution ranking reflects physics performance
- **HEP code is/remains competitive**
  - Alternative solutions become powerful
- **Domain knowledge does help**

## Phase 2

- **Entry barrier was (too) high**
  - Predefined API makes it more restrict
  - Less known/prominent coda lab platform
- **HEP code is currently the best option**
  - Can still be improved (is ongoing)
- **Alternative solutions lack ~1 order of magnitude in CPU behind**
  - Worth watching, and certainly not everything was tried

# Spin-offs & Remains



Dataset used in many R&D projects

From this workshop only

A virtual detector

Information Flow

HEP.TrkX

- Checking edge score after each step of graph network.
- Effective output of the model is in step 8.
- Full track hit assignment learned in last stages of the model.
- Tracklets learned in intermediate stages.

query restricted to volume : 7, avg quality: 6.18

Results

dataset size: ~20%  
1,637 particles, 11,030 hits

plotting error: too many doublets 392529

392,529 doublets  
p=0.26%, r=99.15%

57.3s  
build QUBO

2,546 doublets  
(2,964 triplets)  
QUBO size: 14,345

17.1s  
sample QUBO  
running on CPU

1,512 doublets  
p=99.13%, r=97.06%

trackml score 97.55%

Slide from [ [IPA, Oct 19-25, 2019](#) ]

# Spin-offs & Remains



Dataset used in many R&D projects



Detector being evolved into an Open Data Detector project

From this workshop only

A virtual detector

Information Flow

HEP.TrkX

### Spin-off Sneak Preview

TrackML Pixel detector

OpenData Pixel detector

Features:

- described in DD4Hep
- realistic material budget
- non-symmetric in azimuthal angle
- full (G4) and fast (ACTS) simulation
- misalignment possibility

... to be released soon!

Checking edge score after each step of graph network. Effective output of the model is in step 8. Full track hit assignment learned in last stages of the model. Tracklets learned in intermediate stages.

1512 doublets  
p=99.13%, r=97.06%  
trackml score 97.55%

Oct 19-25, 2019 ]

CTDWIT 2019 - A. Salzburger | Summary of the Tracking Machine Learning Challenge 71

More information: [ [CHEP2018](#) ] [ [CTD2019](#) ]

# Spin-offs & Remains

From this workshop only



Dataset used in many R&D projects

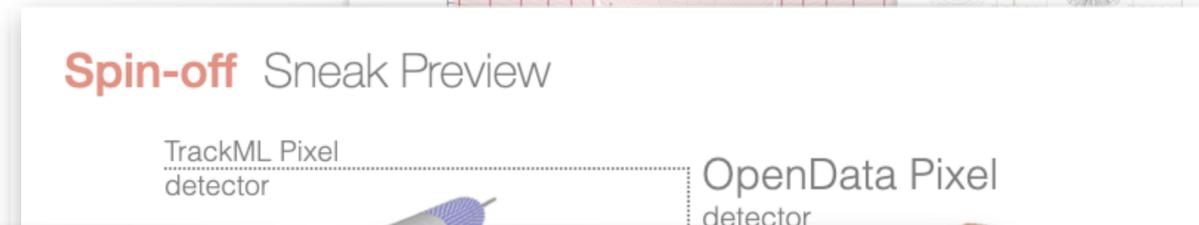
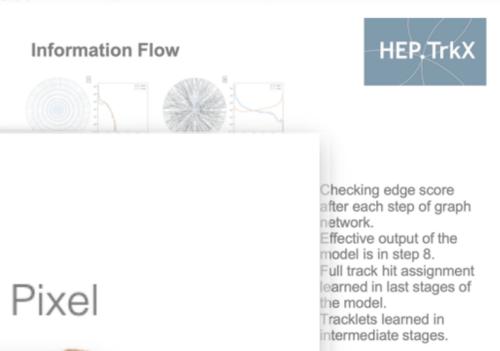


Detector being evolved into an Open Data Detector project



Community

- Competition track at [ [NeurIPS2018](#) ]
- [ [TrackML Grand Finale](#) ], CERN, 2019
- [ [IPA Workshop](#) ], Paris, 2019
- Conference talks, seminars, etc.



### TrackML Challenge: grand finale

1 Jul 2019, 13:30 → 2 Jul 2019, 18:00 Europe/Zurich  
222/R-001 (CERN)

**Description** The HL-LHC will see ATLAS and CMS see proton bunch collisions reaching track multiplicity up to 10.000 charged tracks per event. Algorithms need to be developed to harness the increased combinatorial complexity. To engage the Computer Science community to contribute new ideas, we have organized a Tracking Machine Learning challenge (TrackML). Participants were provided events with 100k 3D points, and are asked to group the points into tracks; they are also given a 100GB training dataset including the ground truth.

The [TrackML challenge](#) has run in two phases.

- The first "Accuracy" phase has run on [Kaggle platform](#) from May to August 2018; algorithms were judged only on a score related to the fraction of correctly assigned hits. The first phase has seen 653 participants, with top performers with innovative approaches, exposed at NeurIPS 2018 Competition workshop and documented in [1904.06778](#)
- The second "Throughput" phase ran Sep 2018 to March 2019 on [Codalab](#), required code submission; algorithms were then ranked by combining accuracy and speed. It has seen astonishingly fast submissions.

This workshop will gather at CERN organizers, top competitors and anyone with interest with fast tracking.

Registration is free.

Non CERN user are welcome (they need to tick the appropriate box in the registration form).

The workshop takes place from Monday 1:30PM to Tuesday 5:30PM. Tuesday morning is reserved for visit of CERN highlights for non HEP participants. In particular, visits of ATLAS or CMS experiments which will be exceptionally open are foreseen. (Visits are fully booked, new registrants will be on the waiting list).

Contact : [trackml-contact@googlegroups.com](mailto:trackml-contact@googlegroups.com)

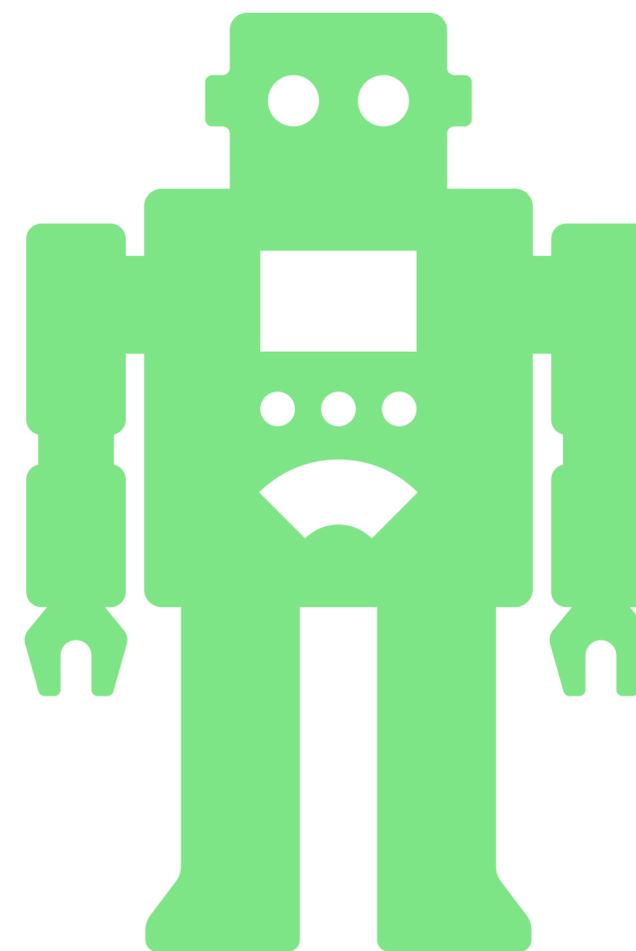


[Oct 19-25, 2019](#) ]

[CTD2019](#) ]

# Conclusions

- **TrackML project started with an idea in 2015**
  - Based on experience of the HiggsML challenge
  - Took some 2-3 years in making, performed through 2018
  - Ran very successfully
- **Huge resonance inside/outside the field**
  - Connected DS/ML and HEP Tracking community
  - Would we have liked to see more ML ? Probably.
- **Spin-offs will remain**
  - Dataset is still at great use in R&D projects
  - Open Data Detector to come for refined development/testing





[trackml.contact@gmail.com](mailto:trackml.contact@gmail.com)



<https://sites.google.com/site/trackmlparticle/>



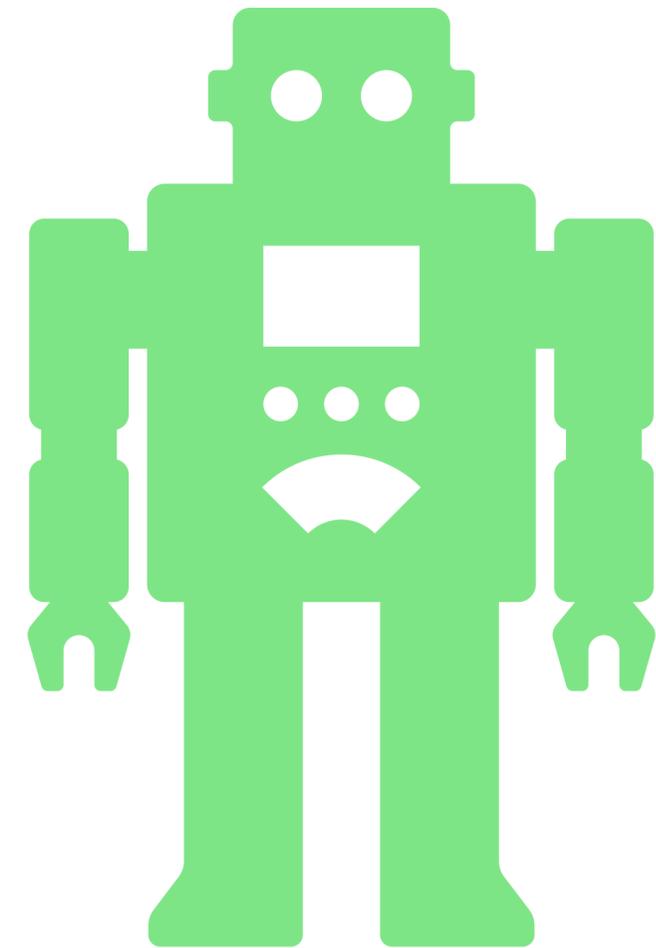
@trackmlhc



<https://www.kaggle.com/c/trackml-particle-identification>



<https://competitions.codalab.org/competitions/20112>



# Backup Material

# Phase 1 Organisers Prize



## • Algorithm outline

- First step is a route data bank building  
*Geometry identifier (module, layer, volume)*  
*used to pre-build route patterns,*  
*route is a sequence of modules*  
  
*assuming training set contains all possible patterns*
- Second step is hit matching  
*searching through all possible routes and check*  
*if you have hits on each module*  
*this defines a track candidate*

