

素粒子実験へのAIの応用研究の最前線

田中純一

東京大学 素粒子物理国際研究センター

2024年6月14日

自己紹介

- 東大の素粒子実験の研究室（相原研の一期生）卒業
 - Belle実験（つくば、KEK）でD論を書きました。
 - 研究室の同期には横山先生（東大理学部物理）、樋口先生（東大IPMU）
 - 同期の東大大学院関係の先生：濱口先生、福嶋先生、矢向先生（結構います…）
物性にも…
- 2002年から素粒子センターでアトラス実験をやっています。



ヒッグスでは終わりたくないですね。

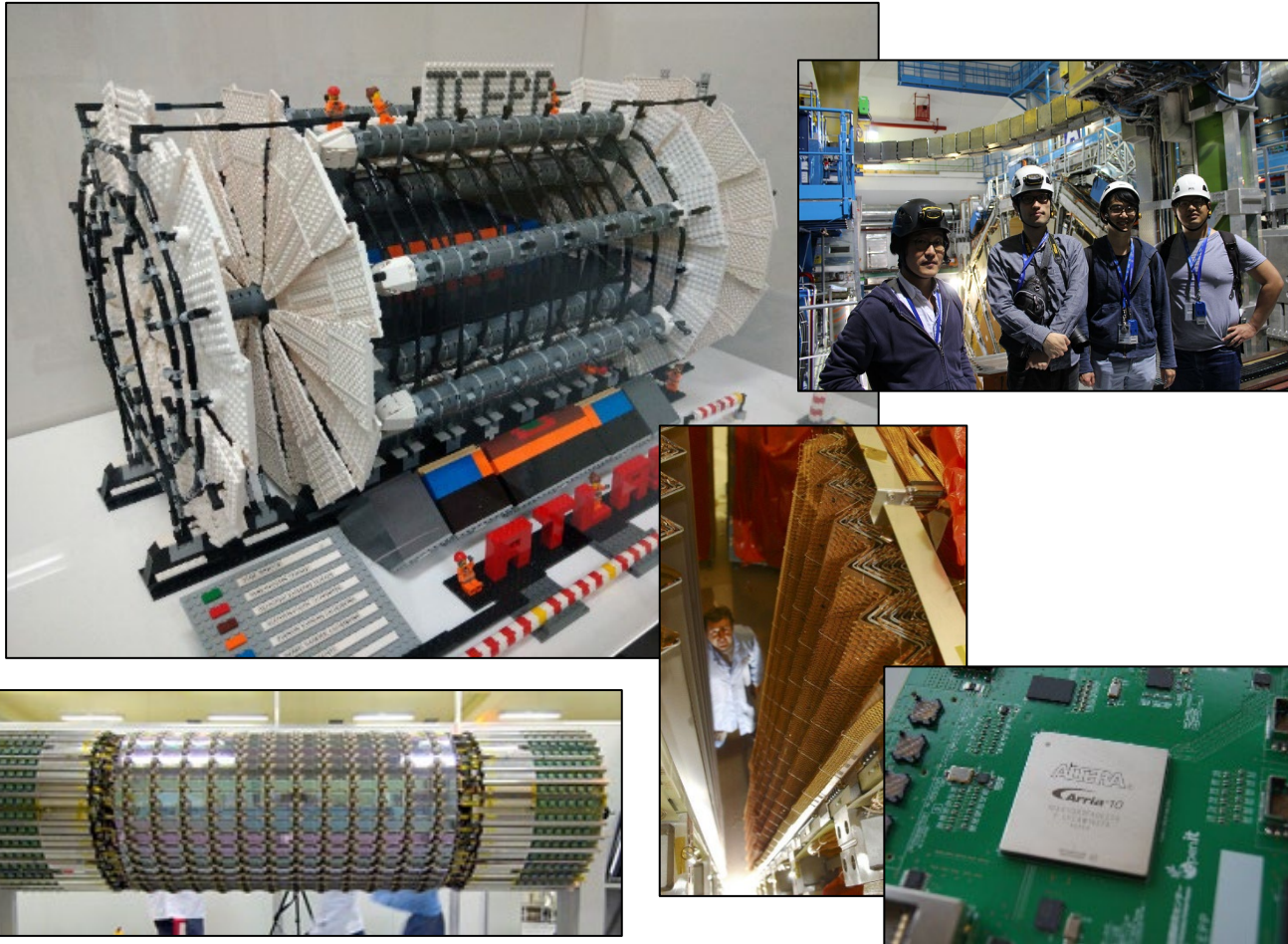
- 計算機センターのお仕事：Grid
- データ解析の準備（2010年まで待った）
- 実際のデータ解析 → 2012年7月 ヒッグス発見
- 2013年3月 EMカロリメータの読み出しデジタルトリガーのR&Dに正式に参加
 - FPGA Firmwareの開発
 - デジタルトリガーのオペレーション
- 2018年4月 Tokyo Tier2, ビッグデータを使った計算機科学
 - 計算機センター（グリッド、クラウド、スパコン）
 - 深層学習、量子コンピュータ
 - Beyond AIなどにも参画していた

田中研：アトラス募集。もちろん、アトラス以外もOKです！

素粒子「実験」の研究

みなさんにはたくさんの選択肢があります！

ハードウェア開発・オペレーションの仕事



計算機上での仕事

検出器データ

再構成

再構成データ

粒子識別

解析データ

解析・検定

学術データ

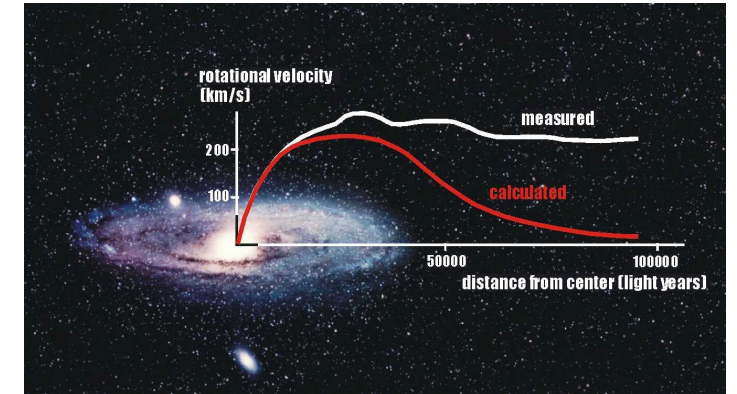
今日は人工知能を使った研究の紹介！

素粒子物理・素粒子実験

Understand our universe, the origin of everything
 → how to interact?, what is matter?, what is force?, ...

Find **ultimate equation(s)** to explain our universe!

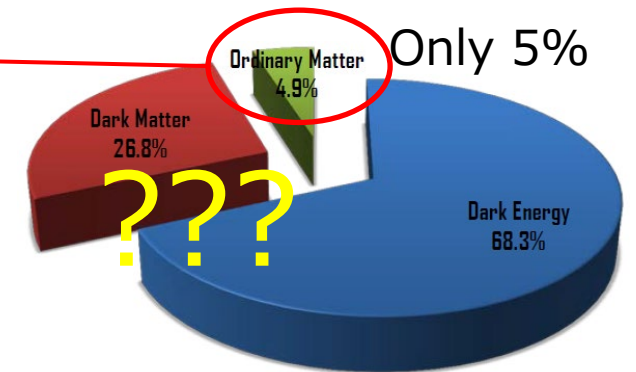
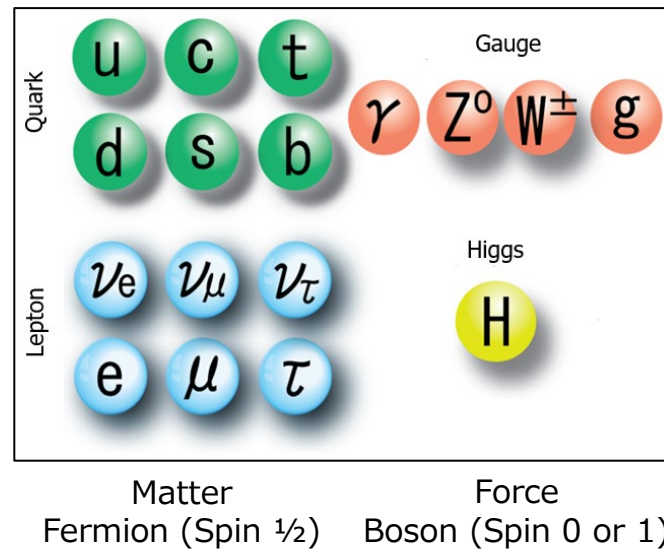
見えている物質で計算した速度より早く動いている。



素粒子標準模型

$$\mathcal{L} = -\frac{1}{4} F_{\mu\nu} F^{\mu\nu} + i\bar{\psi}\not{D}\psi + h.c. + \chi_i y_{ij} \chi_j \phi + h.c. + |D_\mu \phi|^2 - V(\phi)$$

17個の素粒子



(個人的に) 達成したいこと: 今の標準模型を超えたものを見つけたい. 暗黒物質候補とSUSYの証拠

世界最高エネルギーの**LHC加速器**は宇宙初期(ビッグバンの 10^{-12} 秒後の世界)の状態を人工的に作り出すことができる!

Swiss Geneva

CERN



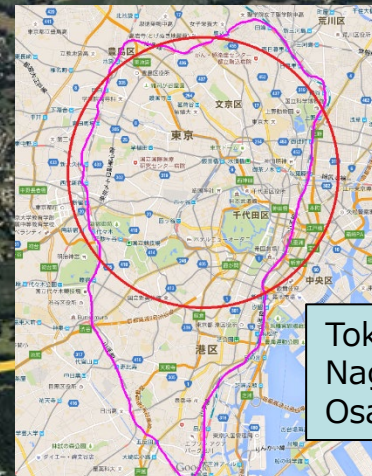
CERN Meyrin

ATLAS

SPS 7 km

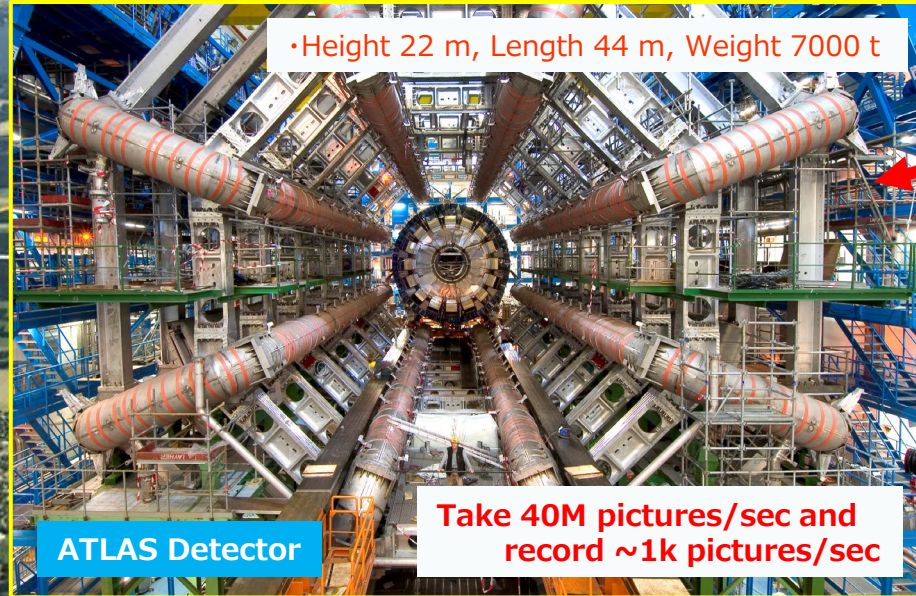
ALICE

LHC tunnel 27km circumference
100m underground



Tokyo	Yamanote-line	34.5km
Nagoya	Meijo-line	26.4km
Osaka	Kanjo-line	21.7km

•Height 22 m, Length 44 m, Weight 7000 t



ATLAS Detector

Take 40M pictures/sec and record ~1k pictures/sec

検出器のセンサー数 1億チャンネル

1秒間に最大1000枚の写真



SUISSE

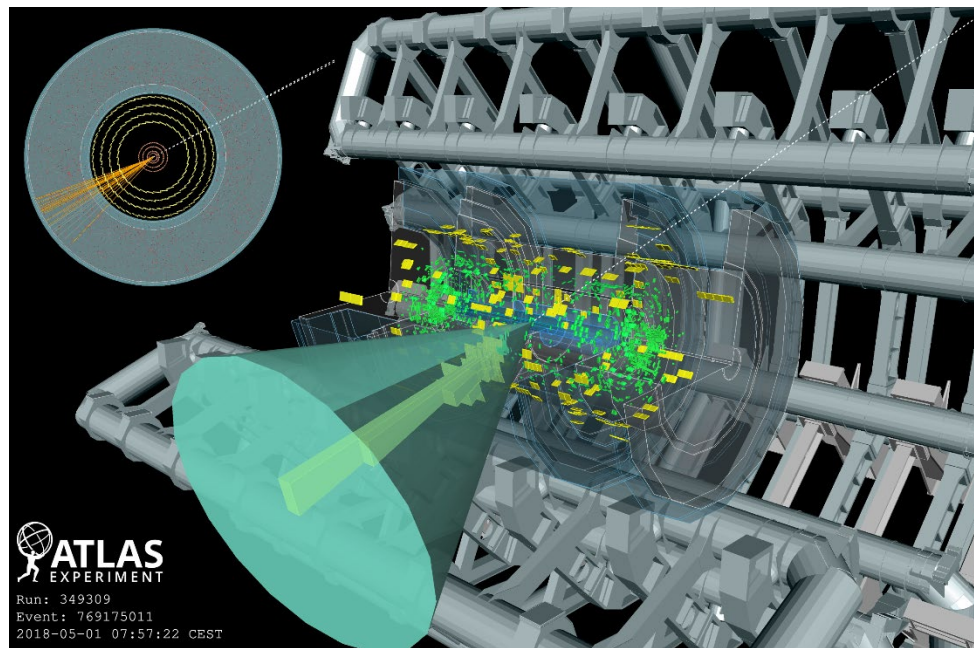
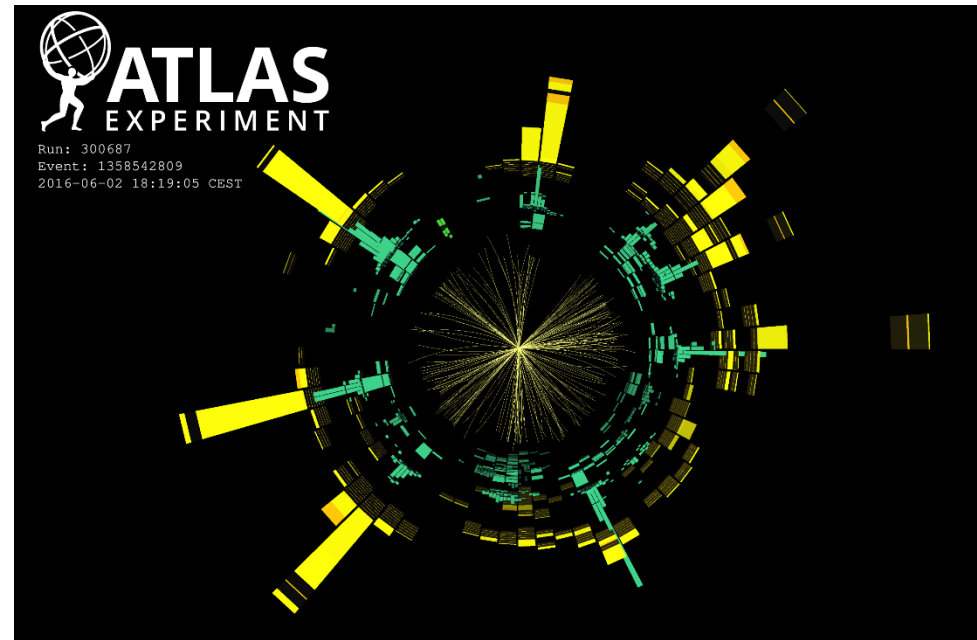
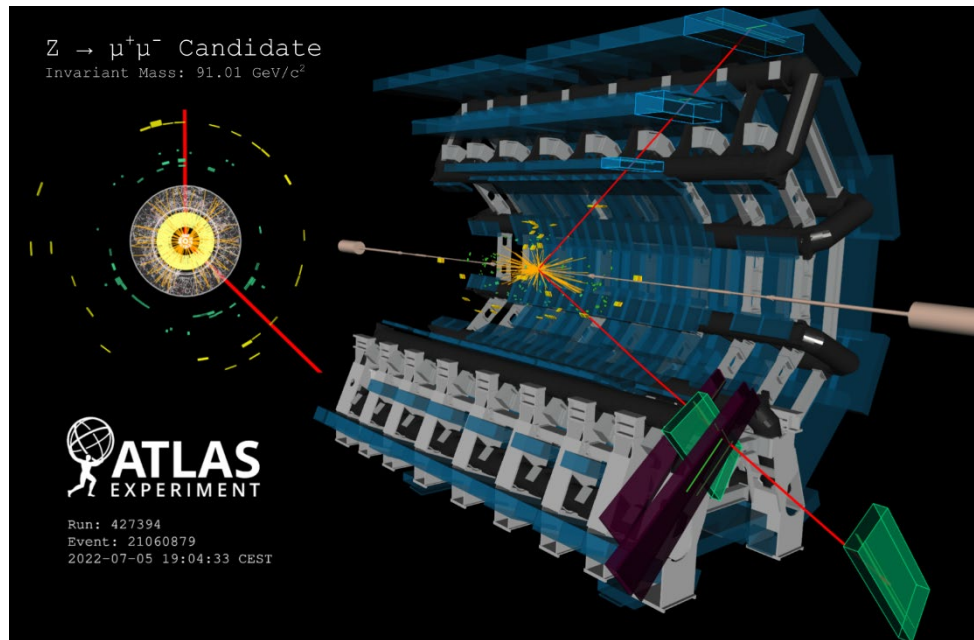
FRANCE



CMS

Large Hadron Collider (LHC) LHC 27 km





実験データの事象(写真!)の一例

ChatGPT3.5: 約45TB

2015-2018年 Run2実験:約2,000,000,000事象
(約20ペタバイト)

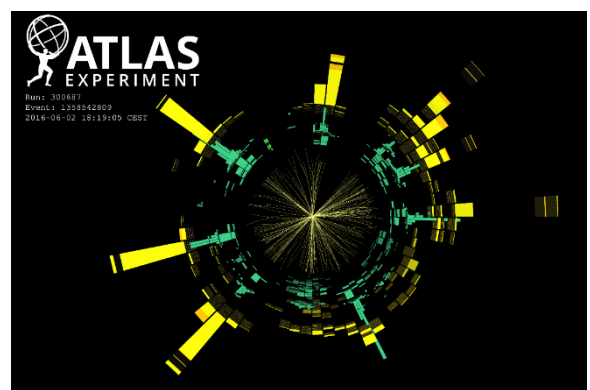
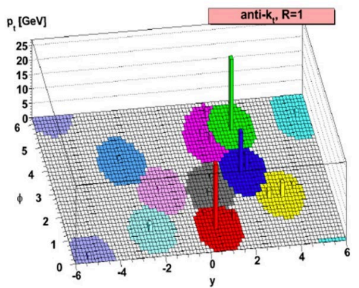
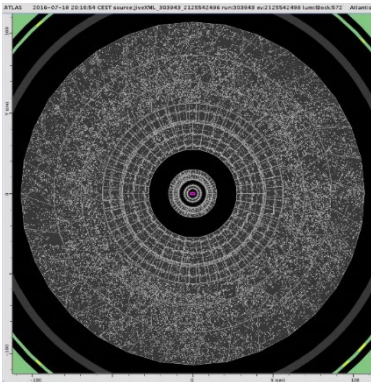
2022-2025年 Run3実験(進行中)
2040年ごろまでに20倍になる計画

シミュレーションデータは数十倍

我々の研究:

実験データを再現したシミュレーションデータとの比較
(シミュレーションデータを駆使した比較)

素粒子実験：データ解析ワークフロー



検出器
データ

- 「再構成」
- 飛跡検出器
 - カロリメータ検出器
 - ミュー粒子検出器

再構成
データ

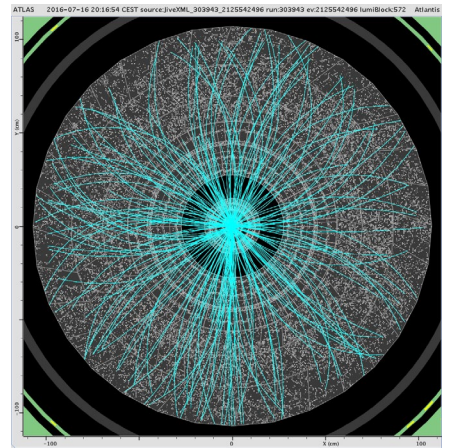
- 「粒子識別」
- 電子や光
 - ミュー粒子
 - ジェットの種類

解析
データ

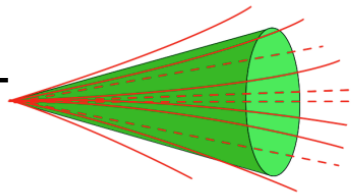
「データ解析」

学术论文 ↓

点
(位置)



曲線や塊：
運動量やエネルギー



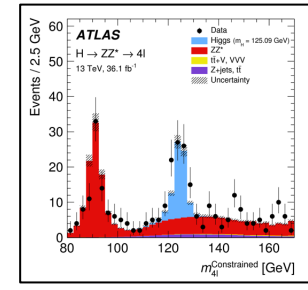
このジェット何？

jet of hadrons

粒子名が付いた
運動量やエネルギー

測定値(E, **p**)を持つ
電子や光子など

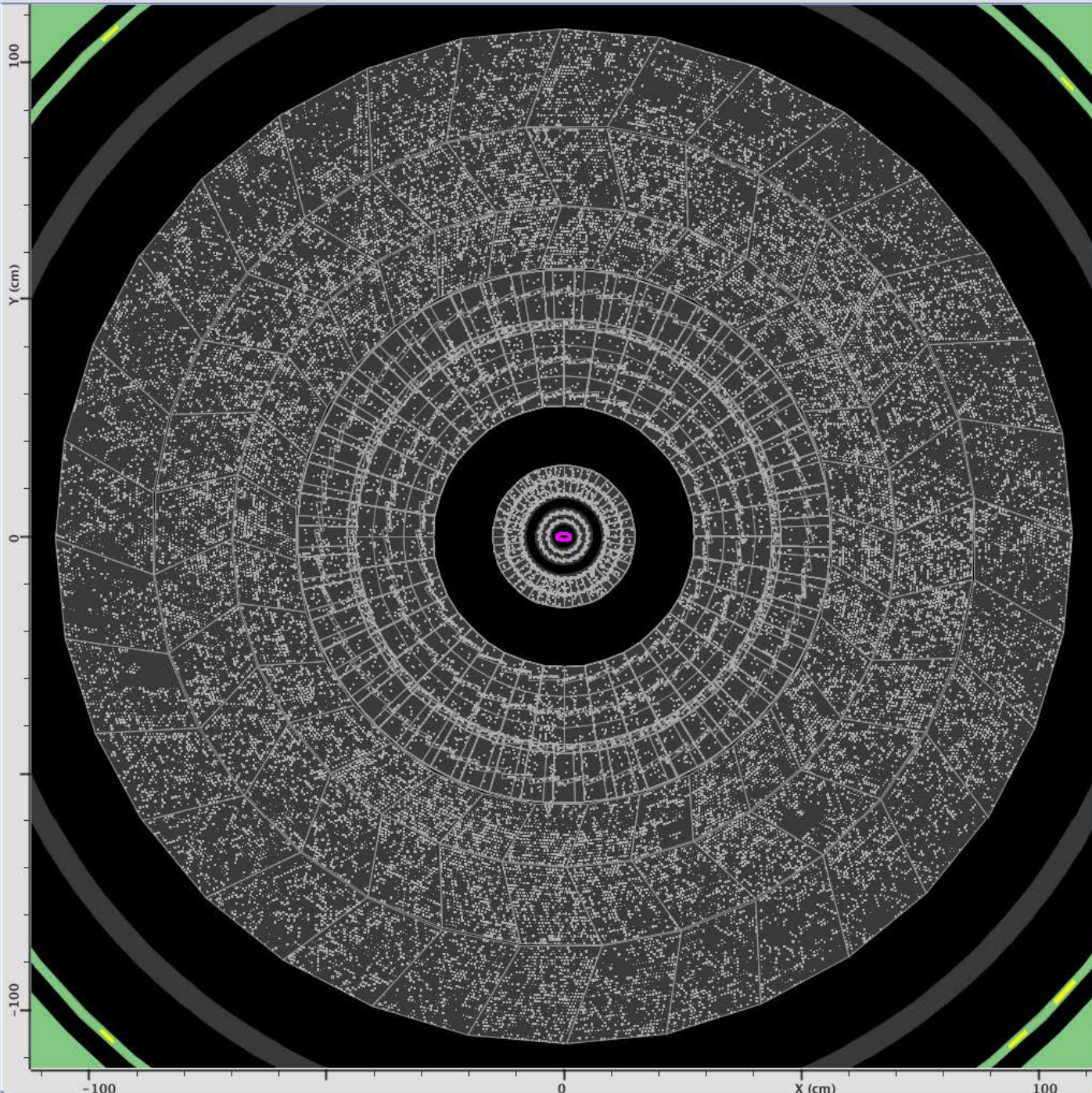
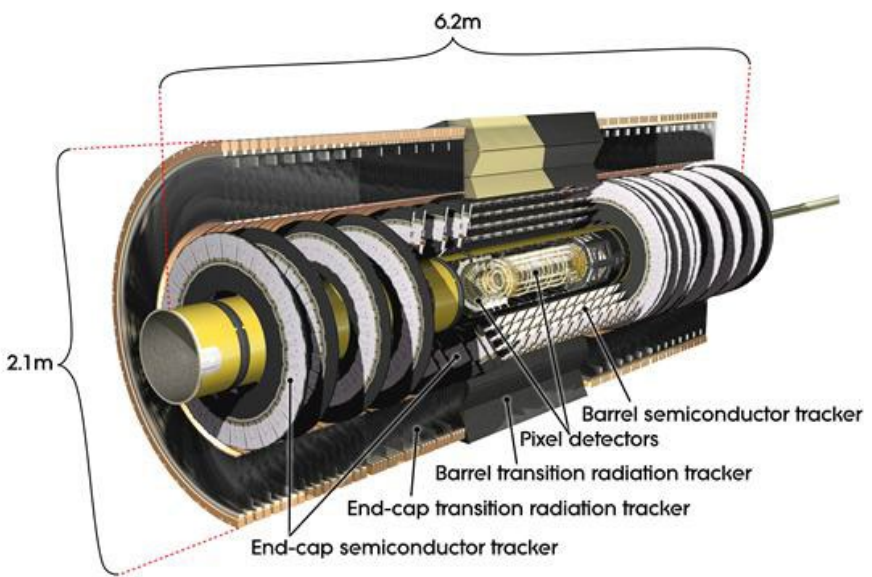
ここまでくれば
シグナルかバックグラウンドか
判断可能になる。



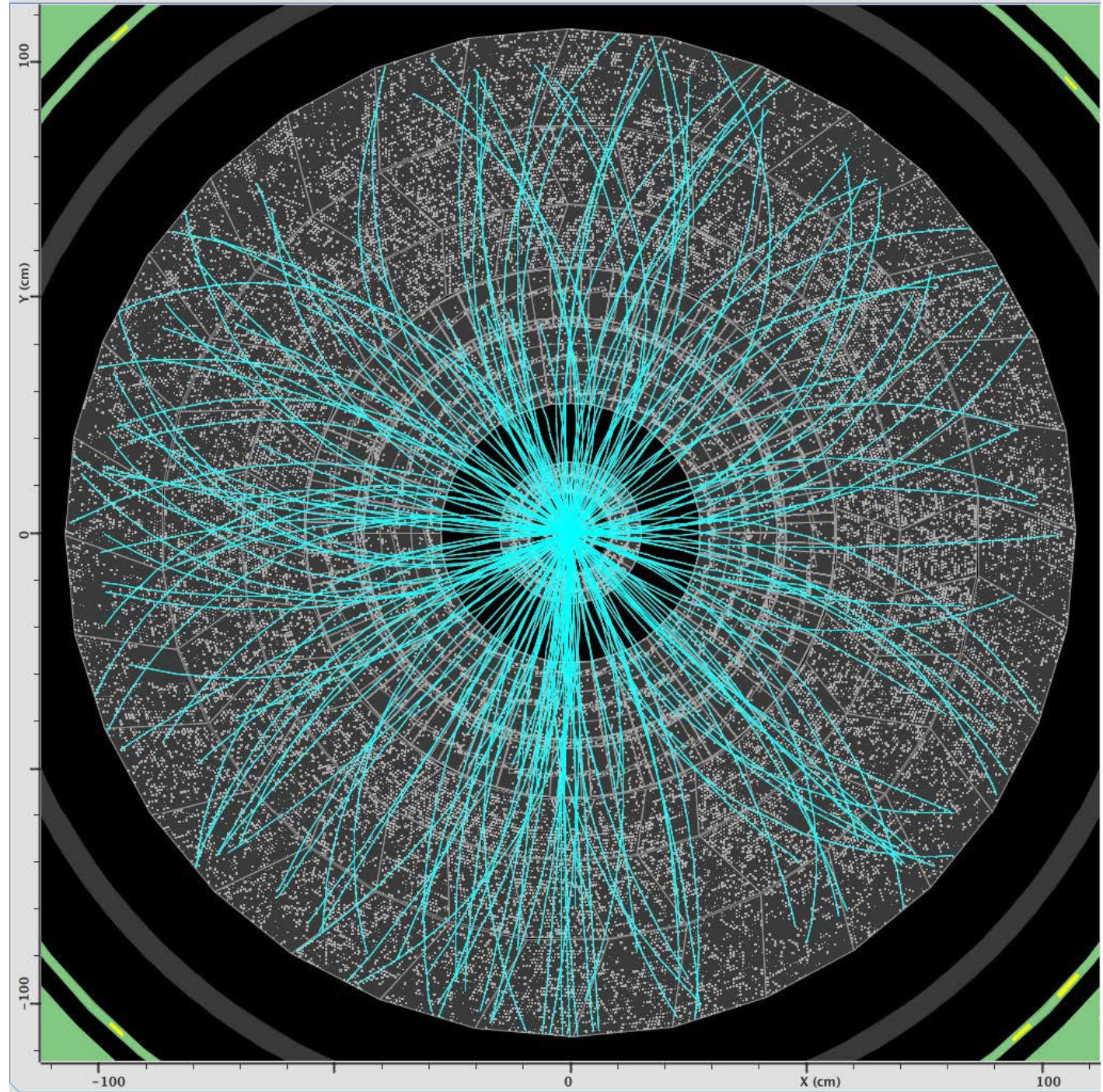
信号事象かどうか？
発見か？

内部飛跡検出器のデータ

- 電荷を持った粒子が「点」を残していく.
- 磁場のため、粒子は曲がる.



- 人工知能の出番！
- 画像系の機械学習？
 - 物理法則は？

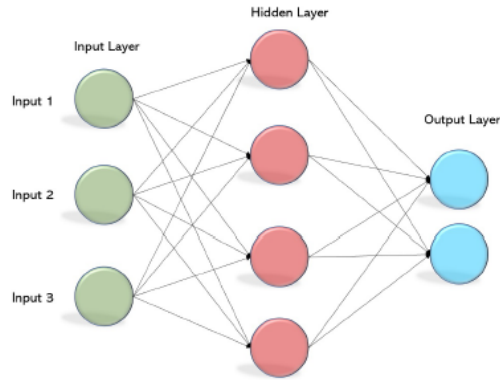


素粒子実験と人工知能

- データ解析の改善
 - オールドフックスなアプローチ：Non-ML・DLからML・DLへ@データ解析ワークフロー
→ これが主流！
 - アグレッシブ：素粒子実験データの基盤モデル（Foundation model）
フレームワークの開発、たとえば、全体最適化（微分接続）
 - よりアグレッシブ：より“intelligent”/“creative”な部分を
 - コード開発
 - データ解析：背景事象の見積もり方法の開発
 - 論文執筆
 - 実験計画・実験設計
- etc.

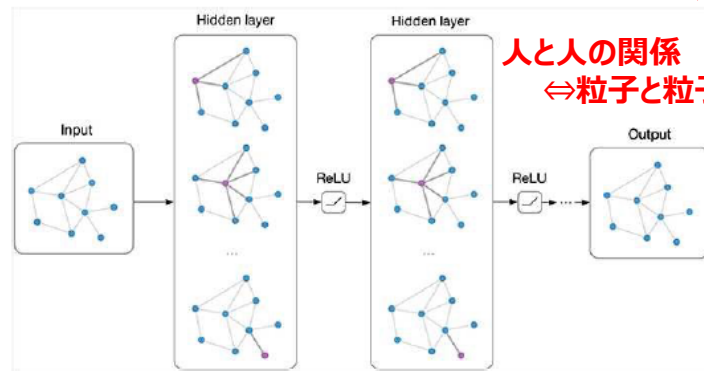
「人工知能」の分野では日々新しいアイデアが生まれている → 素粒子物理への応用

Multi-layer perceptron (MLP)



From: <https://becominghuman.ai/multi-layer-perceptron-mlp-models-on-real-world-banking-data-f6dd3d7e998f>

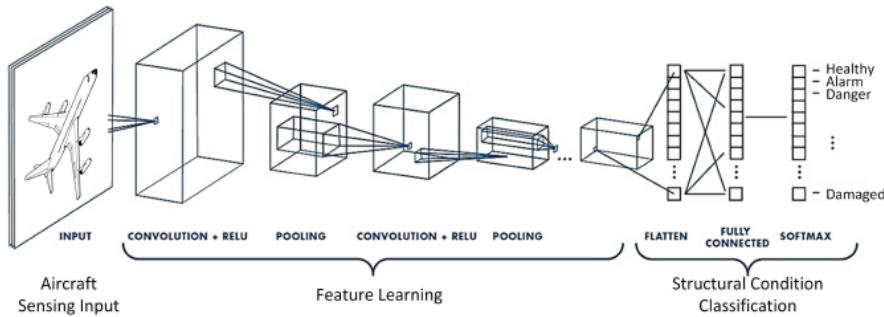
Graph neural networks (GNNs)



人と人の関係
⇔ 粒子と粒子の関係

From: https://theaisummer.com/Graph_Neural_Networks/
 “Scaling law”に従う
 (パラメータ数、データセットサイズ、予算に対してべき乗則)

Convolutional neural networks (CNNs)



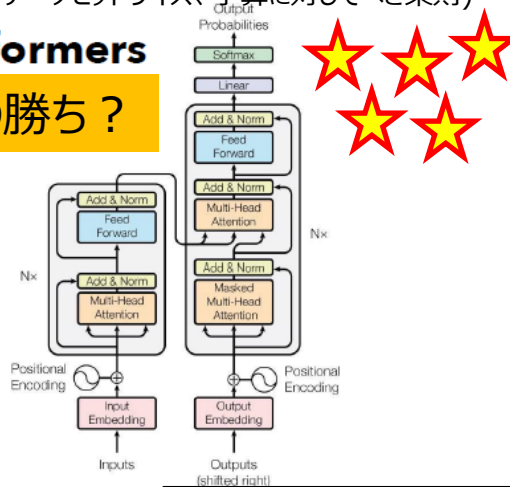
From: [Tabian et al., “A Convolutional Neural Network for Impact Detection and Characterization of Complex Composite Structures,” Sensors 19(22), 2019]

AlexNet (2012) by Hinton’s team

From: <https://pytorch.org/tutorials>

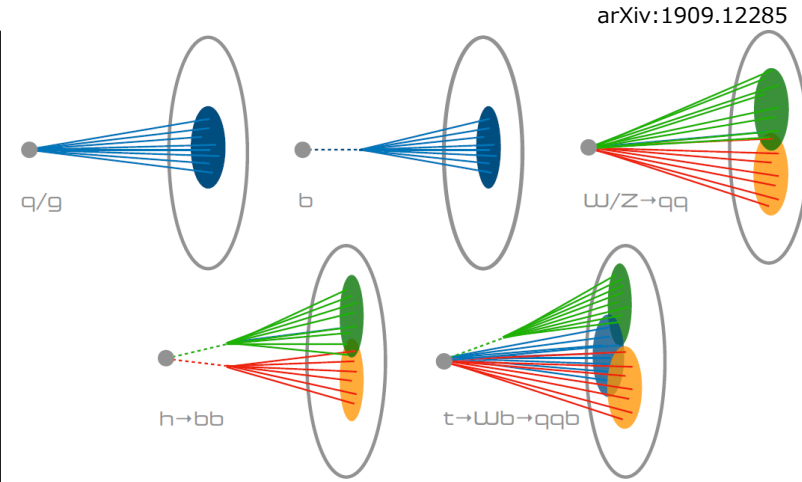
Transformers

金持ちの勝ち?



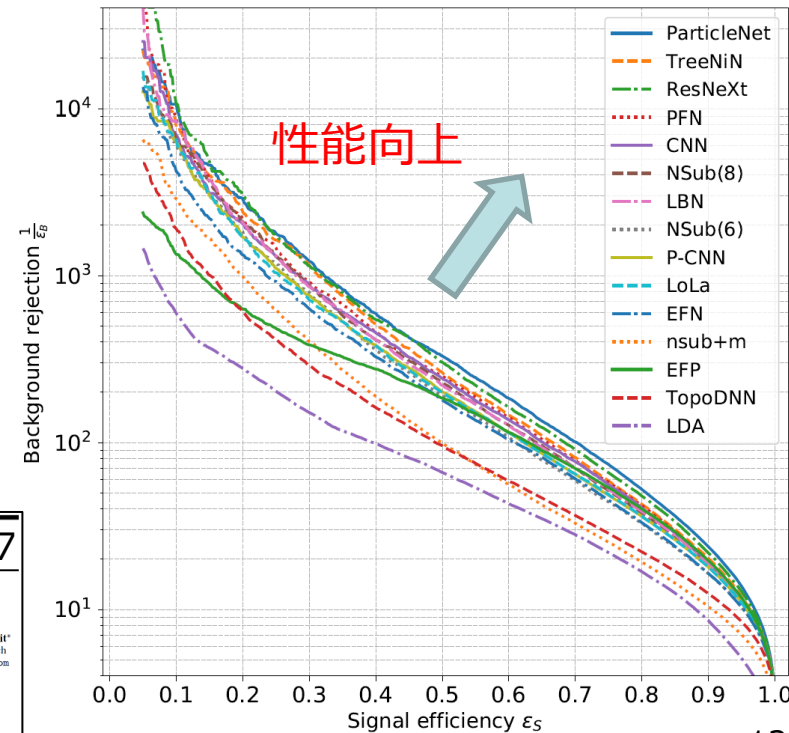
Attention Is All You Need 2017

Ashish Vaswani* Google Brain avaswani@google.com
 Noam Shazeer* Google Brain noam@google.com
 Niki Parmar* Google Research nikip@google.com
 Jakob Uszkoreit* Google Research usz@google.com
 Llion Jones* Google Research llion@google.com
 Aidan N. Gomez*¹ University of Toronto aidan@ca.toronto.edu
 Lukasz Kaiser* Google Brain lukaszkaiser@google.com
 Illia Polosukhin*¹ illia.polosukhin@gmail.com



arXiv:1909.12285

トップクォーク「ジェット」の識別性能



性能向上

arXiv:1902.09914

阪大IDS : 中島悠太先生のスライドから

2012年「ヒッグス発見」と「AIブーム再燃」
 2022年「生成AI元年」

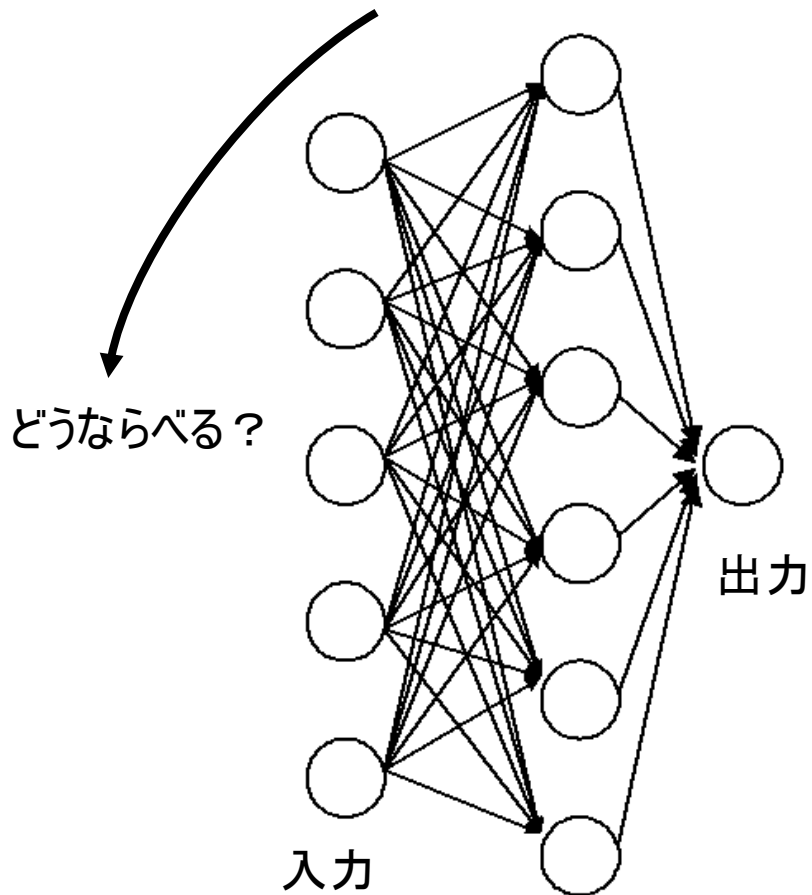
素粒子実験データとGraph Neural Network

荷電粒子の特徴: 運動量と位置 (インパクトパラメータ)
クラスターの特徴: 運動量 (エネルギーとその場所) etc.

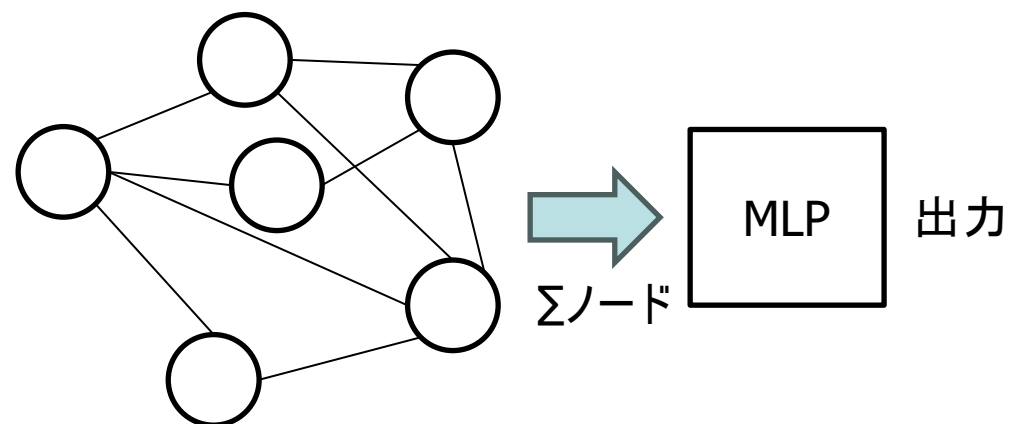
荷電粒子の数
クラスターの数

特徴量の数は固定

事象ごとに異なる



グラフニューラルネットワーク (GNN)



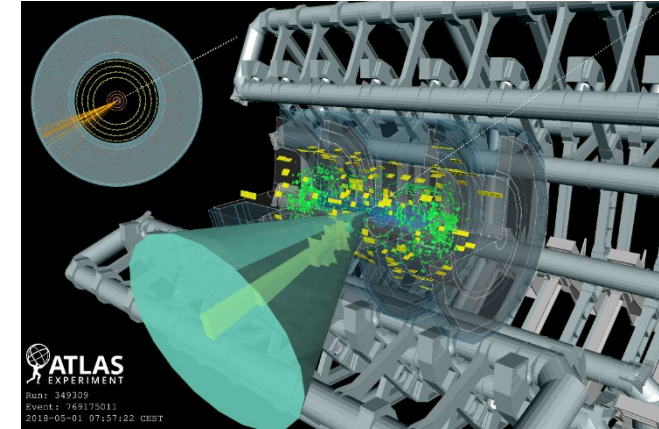
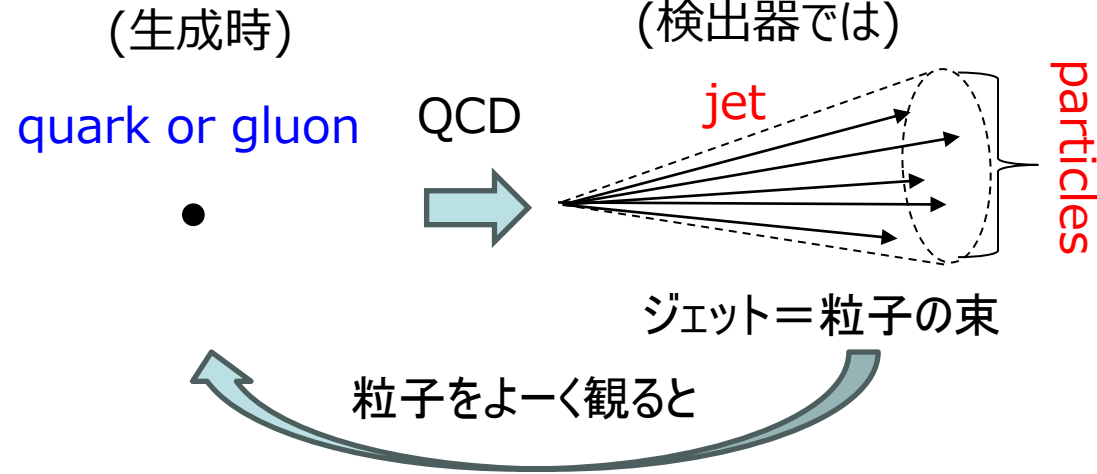
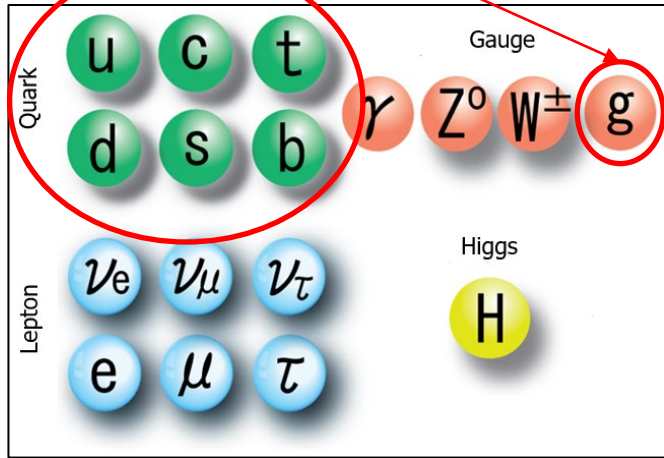
- 各ノードに荷電粒子やクラスター (の表現) を割り当てる。
- ノードやエッジがMLP

MLP = Multi layer perceptron (多層パーセプトロン)

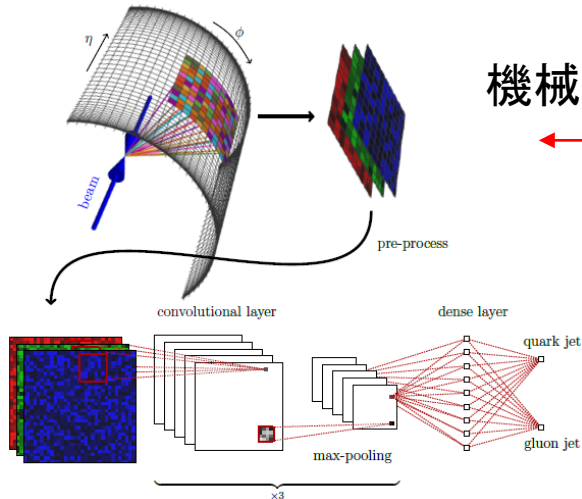
GNNは素粒子データと相性はいい

この人たちは単独では観測できない(QCD理論)!

ジェットの識別

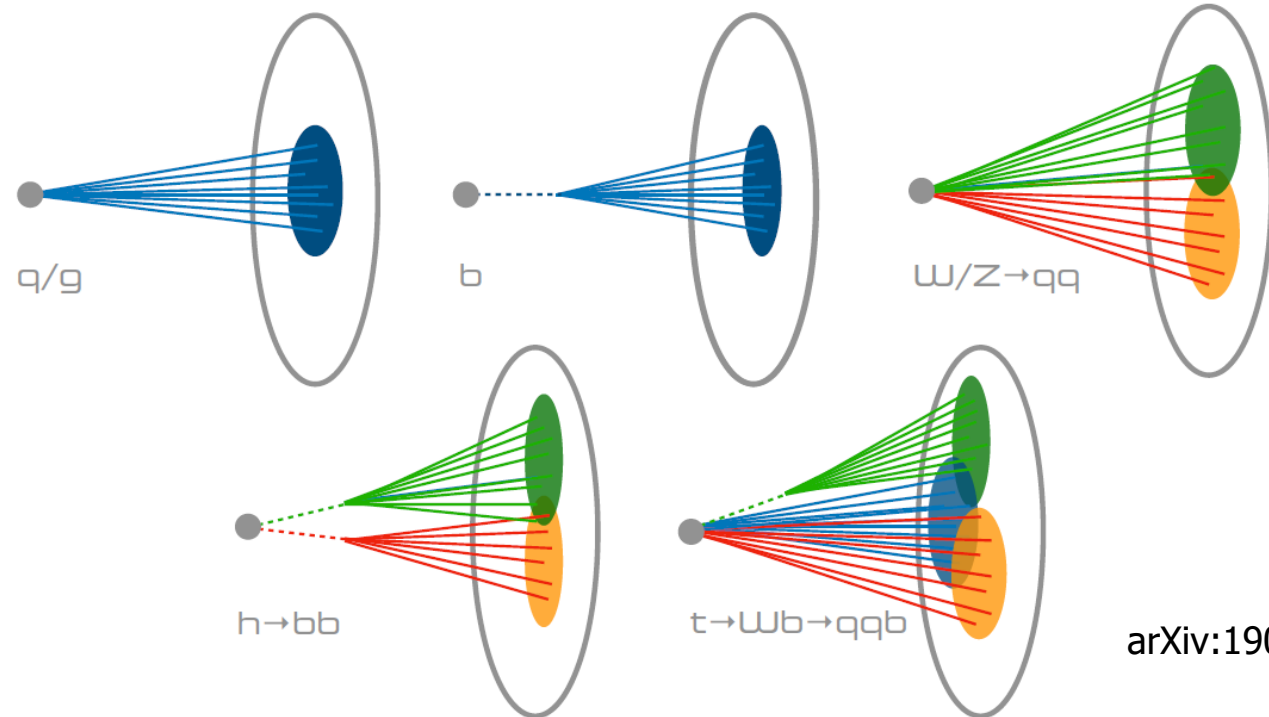


歴史的にはCNNから



arXiv:1612.01551

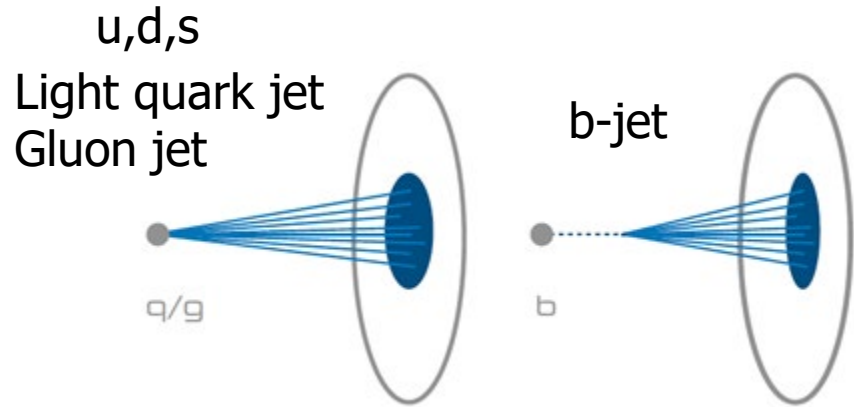
機械学習化



arXiv:1909.12285

Graph Neural Networkを用いたb-tagging

H→bb(ヒッグス粒子の崩壊)のような研究: b-jetの識別が鍵



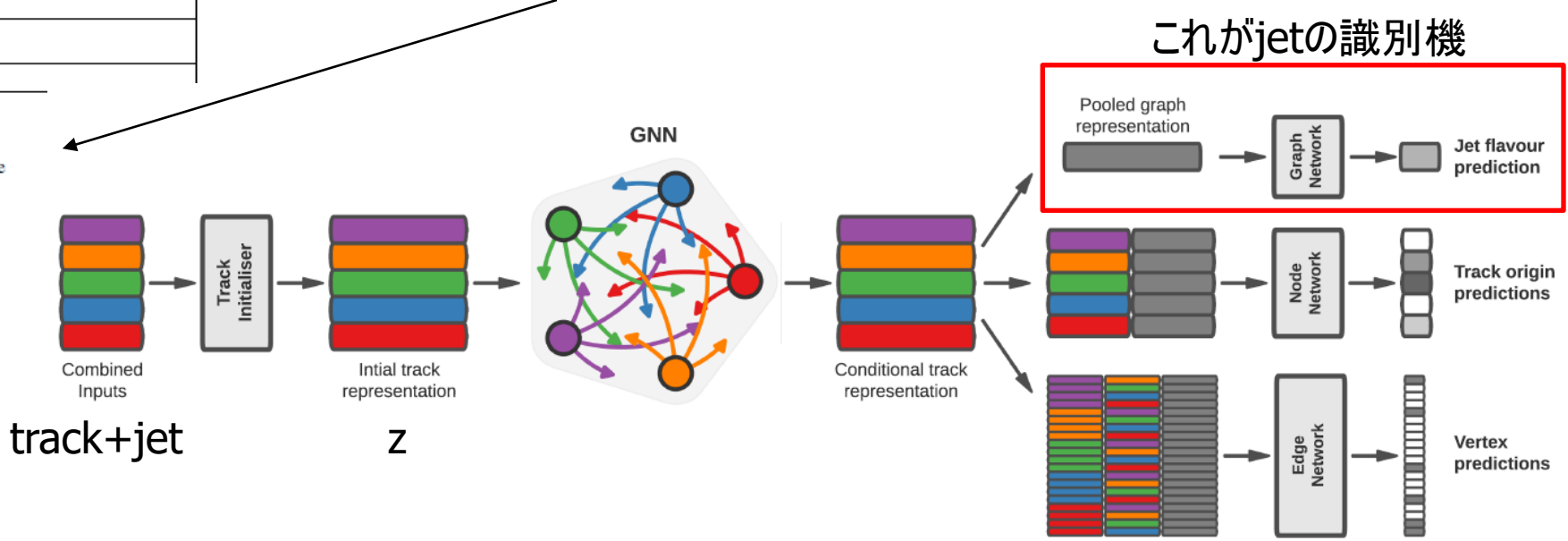
b-jet: Bハドロン崩壊 $B \rightarrow Dlv, D \rightarrow Klv$

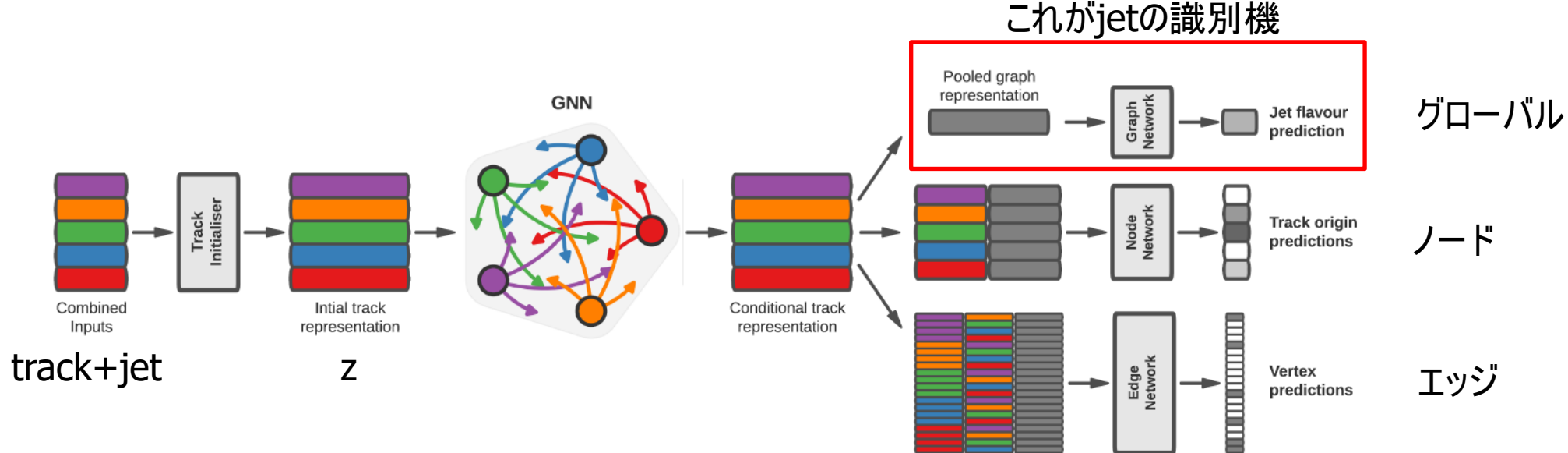
BやDは比較的長寿命 → 検出器で識別可能

従来の方法: 長寿命という性質を抽出する複数のアルゴリズムの開発
それらの出力を機械学習のインプットに。

最新の方法: 粒子(track)を「ノード」として表現したグラフを使った深層学習
入力変数が多いが、どちらかといえば、かなり生データに近い。

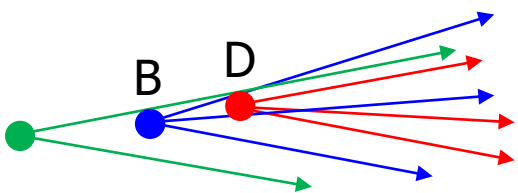
Jet Input	Description
p_T	Jet transverse momentum
η	Signed jet pseudorapidity
Track Input	Description
q/p	Track charge divided by momentum (measure of curvature)
$d\eta$	Pseudorapidity of the track, relative to the jet η
$d\phi$	Azimuthal angle of the track, relative to the jet ϕ
d_0	Closest distance from the track to the PV in the longitudinal plane
$z_0 \sin \theta$	Closest distance from the track to the PV in the transverse plane
$\sigma(q/p)$	Uncertainty on q/p
$\sigma(\theta)$	Uncertainty on track polar angle θ
$\sigma(\phi)$	Uncertainty on track azimuthal angle ϕ
$s(d_0)$	Lifetime signed transverse IP significance
$s(z_0)$	Lifetime signed longitudinal IP significance
nPixHits	Number of pixel hits
nSCTHits	Number of SCT hits
nIBLHits	Number of IBL hits
nBLHits	Number of B-layer hits
nIBLShared	Number of shared IBL hits
nIBLSplit	Number of split IBL hits
nPixShared	Number of shared pixel hits
nPixSplit	Number of split pixel hits
nSCTShared	Number of shared SCT hits
nPixHoles	Number of pixel holes
nSCTHoles	Number of SCT holes
leptonID	Indicates if track was used in the reconstruction of an electron or muon (only for GN1 Lep)





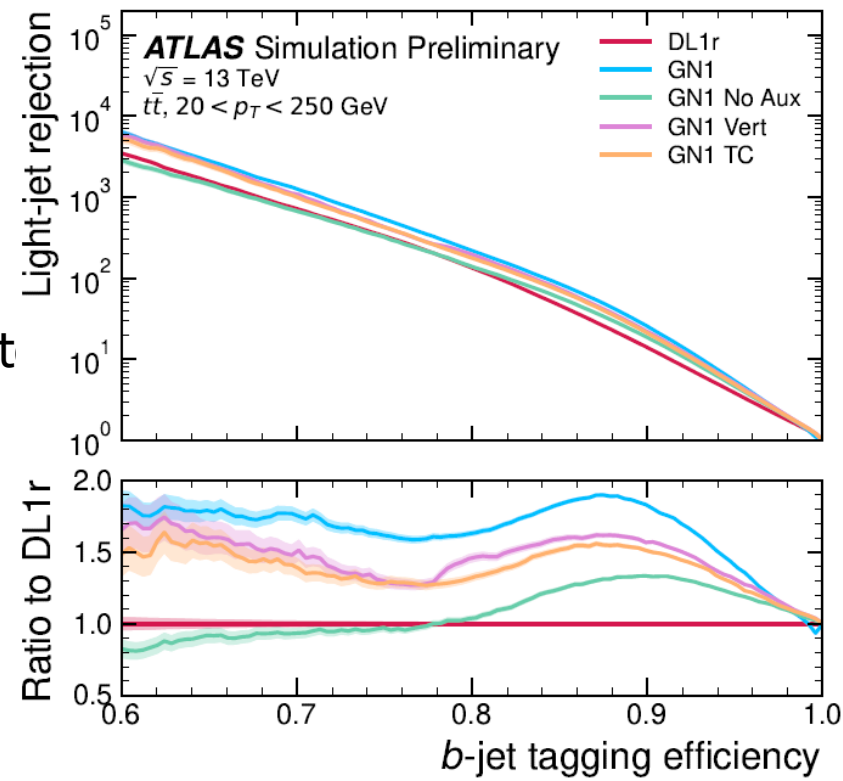
工夫：グラフのノードとエッジ情報も学習に使う！

- 赤と青のトラックがb-jet由来
- 緑のトラックは別の陽子・陽子の相互作用（パイルアップ）

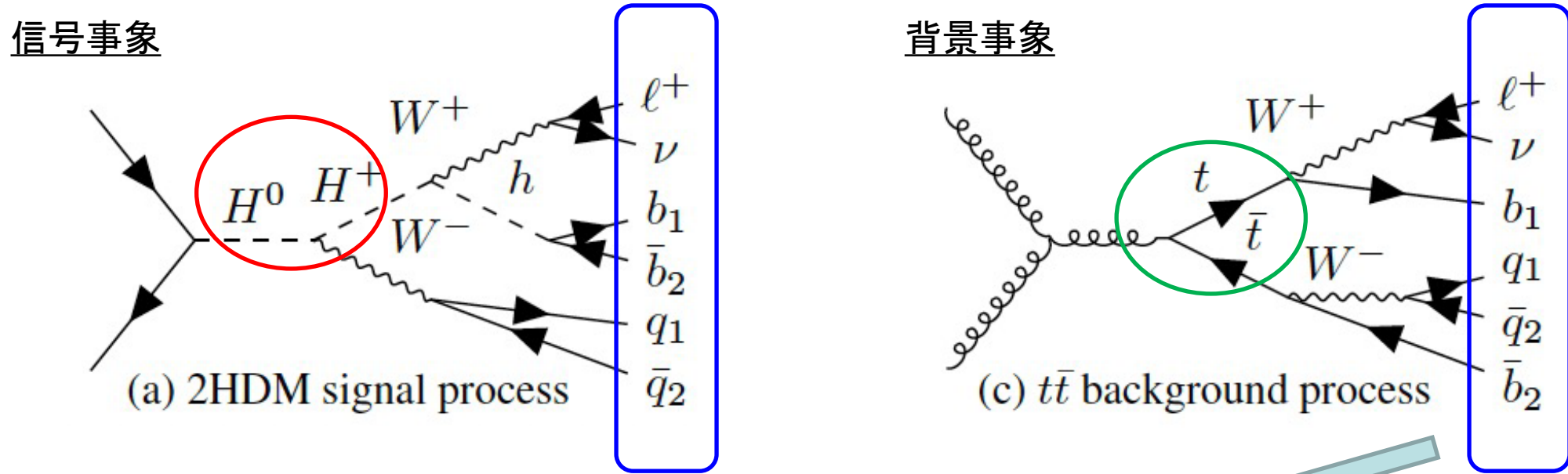


ノード識別：生成点分類（例：赤丸、青丸、緑丸 et
エッジ識別：同じ生成点由来かどうか？

$$L_{\text{total}} = L_{\text{jet}} + \alpha L_{\text{vertex}} + \beta L_{\text{track}}$$



事前知識を機械学習へ



解析データ(検出器データ)でみえるものは**全く同じ**

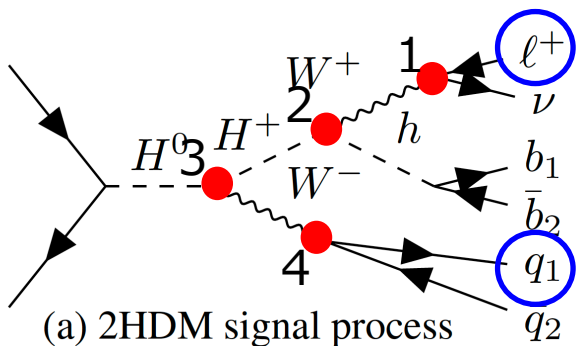
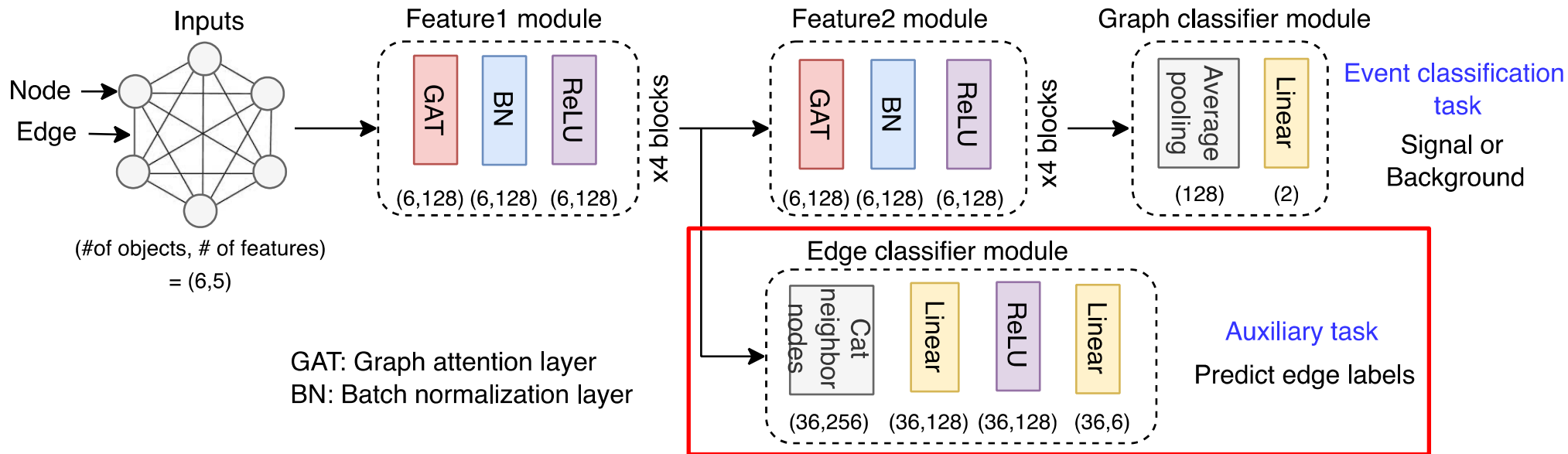
しかし、信号事象と背景事象の**生成過程は異なる**ことは知っている！

→ 信号事象には**途中でHiggs**が存在、背景事象では**途中でtop**が存在などという情報

異なるファインマン図

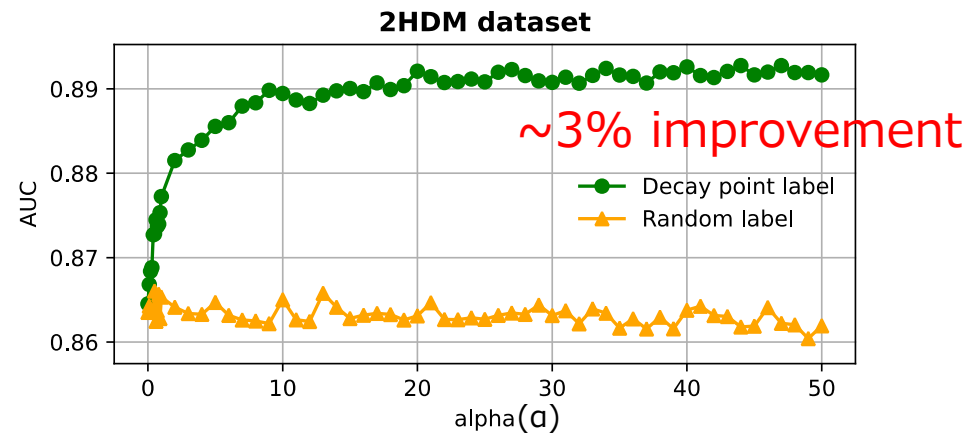
Decay-aware neural network

異なるファインマン図 → 異なるホップ数 (Vertex数)



	l	ν	b ₁	b ₂	q ₁	q ₂
l	0	1	3	3	4	4
ν	1	0	3	3	4	4
b ₁	3	3	0	1	4	4
b ₂	3	3	1	0	4	4
q ₁	4	4	4	4	0	1
q ₂	4	4	4	4	1	0

Label matrix based on decay chains



Loss = L (event classification) + α * L (auxiliary task)

Scaling Laws

(ちょっと大きな話?)

arXiv:2001.08361

Scaling Laws for Neural Language Models

Jared Kaplan * Johns Hopkins University, OpenAI jaredk@jhu.edu		Sam McCandlish* OpenAI sam@openai.com	
Tom Henighan OpenAI henighan@openai.com	Tom B. Brown OpenAI tom@openai.com	Benjamin Chess OpenAI bchess@openai.com	Rewon Child OpenAI rewon@openai.com
Scott Gray OpenAI scott@openai.com	Alec Radford OpenAI alec@openai.com	Jeffrey Wu OpenAI jeffwu@openai.com	Dario Amodei OpenAI damodei@openai.com

Transformerの凄さ(?)

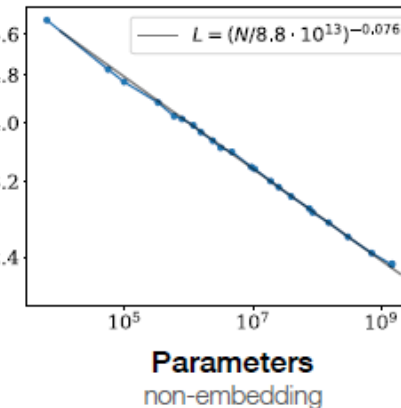
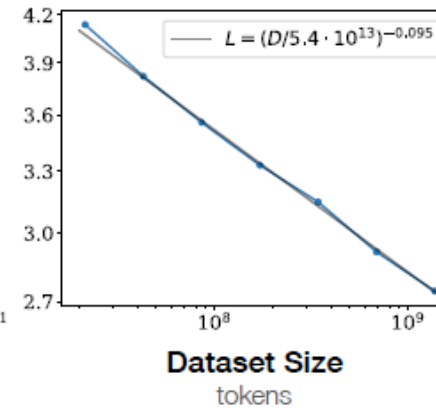
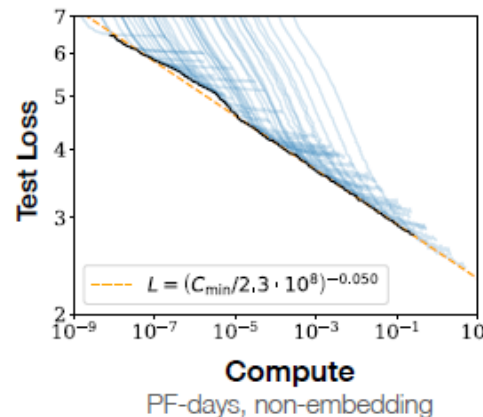
経験則であるが、**スケーリング則**とは

- ネットワークモデルのパラメータ数
- データ数
- 計算量

が上げれば上がるほど、性能が向上する

金持ちの勝ちの法則... ?

ChatGPTなど



Transformer

学習データが膨大なら、Transformerでなんとかなるのでは？

arXiv:1706.03762

$$\text{Attention}(Q, K, V) = \text{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right)V$$

Q: Query
K: Key
V: Value

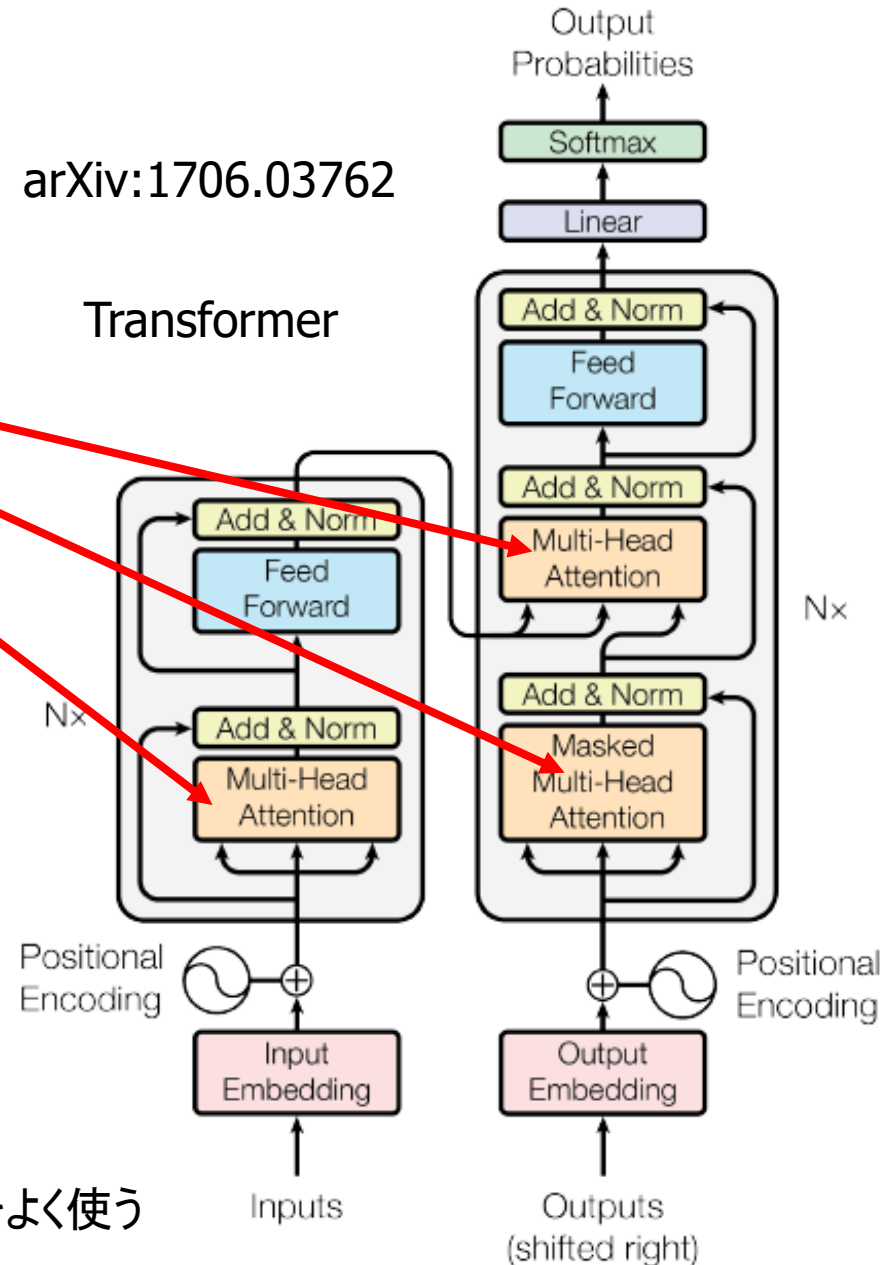
KeyとValueが「メモリー」的な働き

Queryは「問い合わせ」

QK^T で「問い合わせ」と「メモリー」を照らし合わせて、関係性、つまり、「どこに注目すべきか」を抽出。(Attention weight)

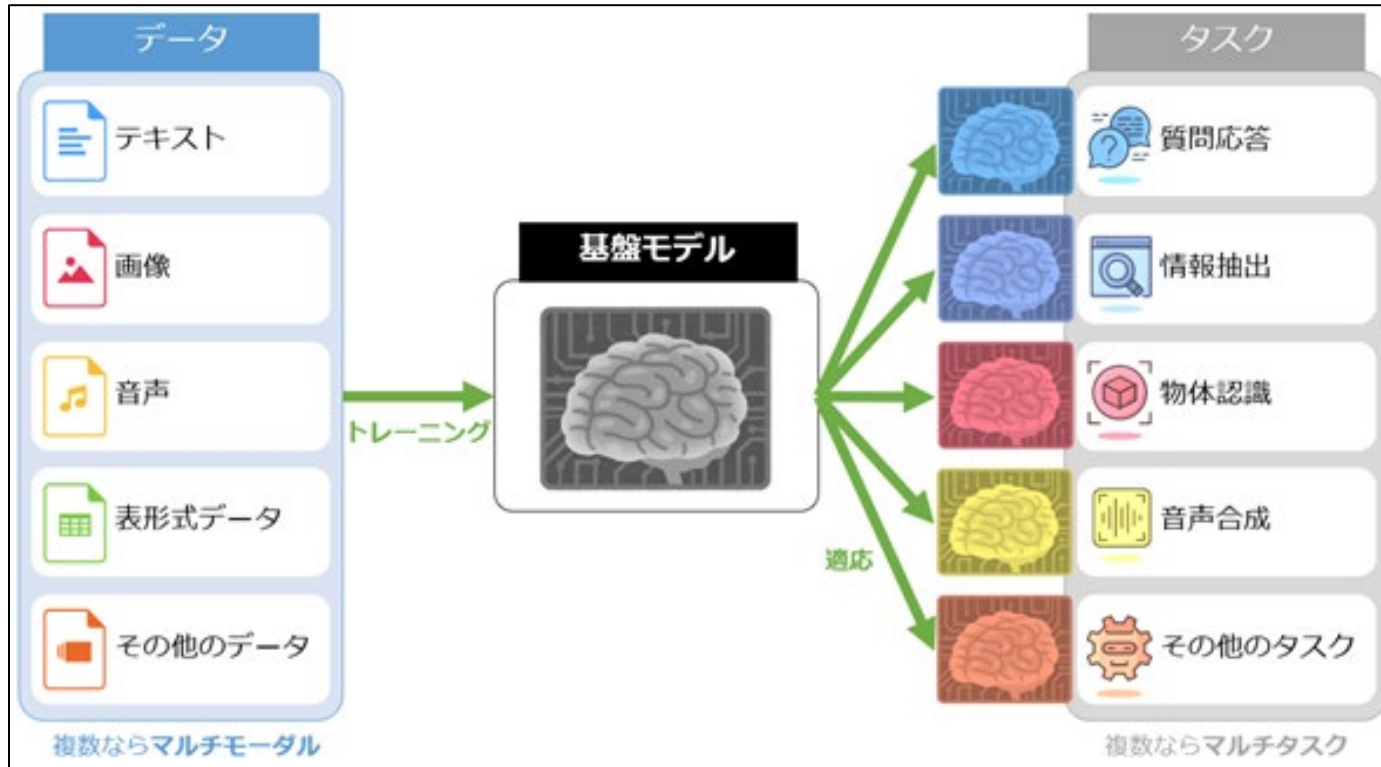
Self-attentionの場合、 $Q=K=V$

素粒子データでは
左側(エンコーダー)だけをよく使う

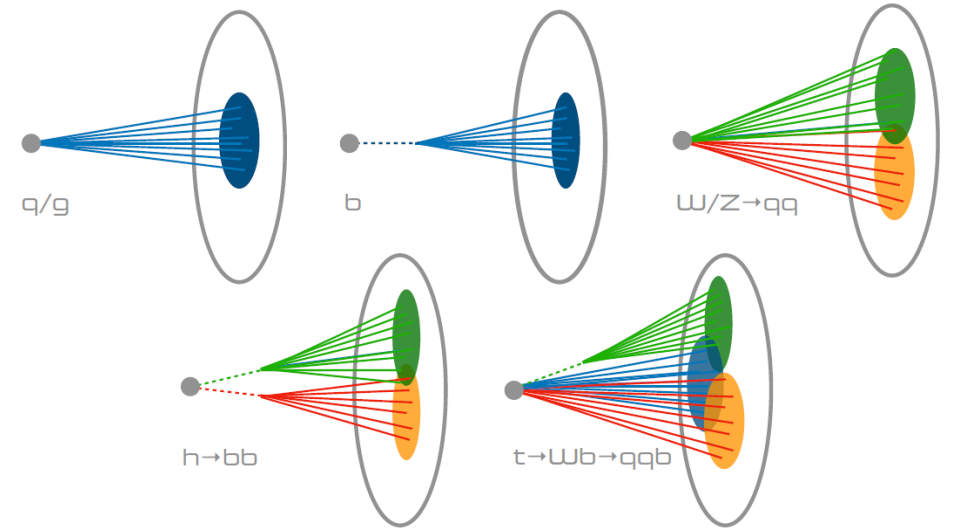


「基盤モデル」：ジェット汎用識別

「基盤モデル」



ジェットもいろいろ: q,g,b,W,Z,H,top etc.



ちょっとテーマが小さいと思うかもしれないが、素粒子データでも「基盤モデル」使えそう

ざっくり言うと「基盤モデル」とは「汎用的で賢いモデル」のこと。

実際に使うには、それをベースに目的に合わせてファインチューニング(再トレーニング、転移学習)をする。

粒子識別等の性能改善（識別問題）

- **b-tagging, c-tagging**, tau-ID, q/g tagger
- W/Z, **Top**, Higgs tagger
- **(Low p_T) electron**, (low p_T) muon etc.
- **信号事象と背景事象の分類**
- **Particle flow** (“粒子”と“測定データ”をマッチさせることでより正確に識別)



AIって：

- Artificial intelligence ?
- Assistant intelligence ?

シンギュラリティってあるの？

シミュレーション

- 検出器シミュレータ (Geant4の代替)
- 物理現象のシミュレータ：**ジェット生成**

ビッグデータから新粒子を発見できるのか？

- **異常検知**

科学するAI・オートメーション化（自動運転） etc.

- 説明責任 (Explainable AI) , 解釈可能 (Interpretability)
- 物理法則、対称性などの抽出
- 素粒子実験の一連の流れを（半）自動化するようなシステム：人と人の連携 → AIとAIの連携
(→**微分可能なプログラム**を作る)

“ChatGPT 3.5”の出現 (2022年11月)

基盤モデル

LLM (Large Language Model)

データ解析の**基盤モデル**のPoC
転移学習は使える？

ATLAS実験のデータ量は膨大

素粒子物理で**人工知能**（**データサイエンス**）を研究しましょう！

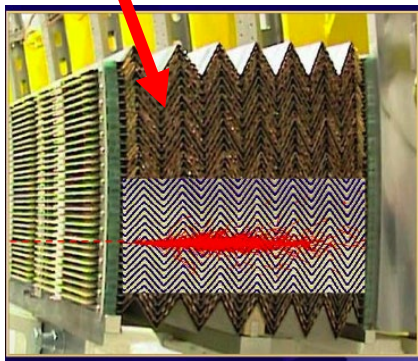
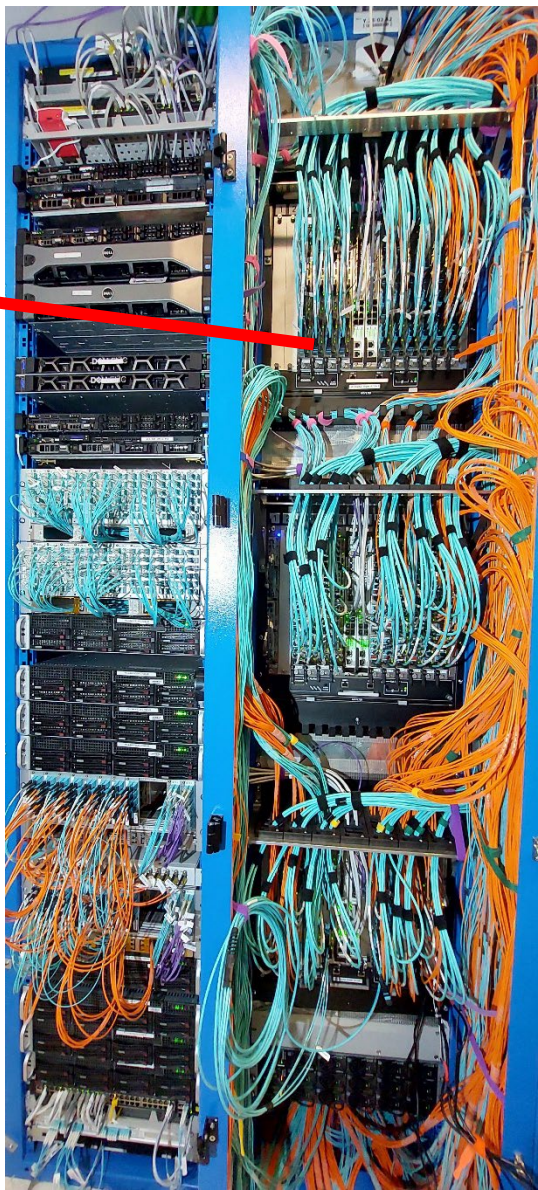
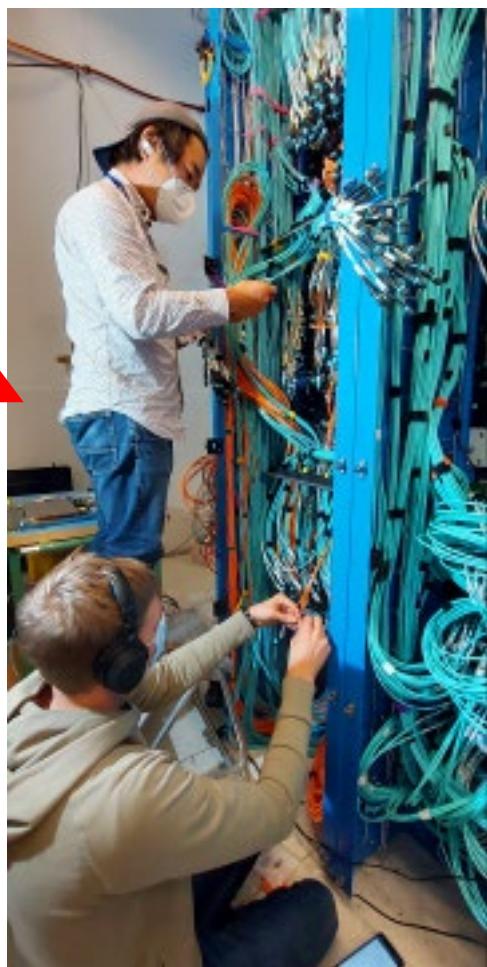
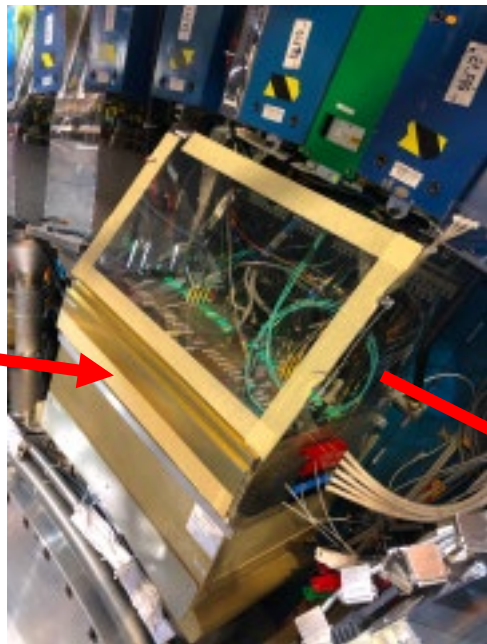
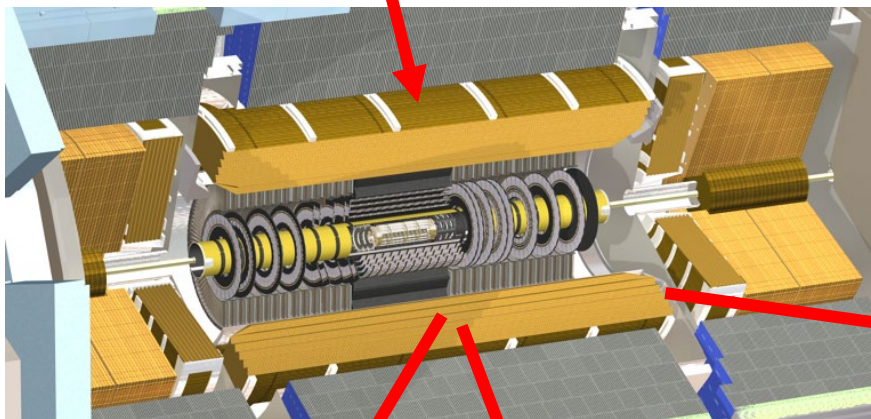
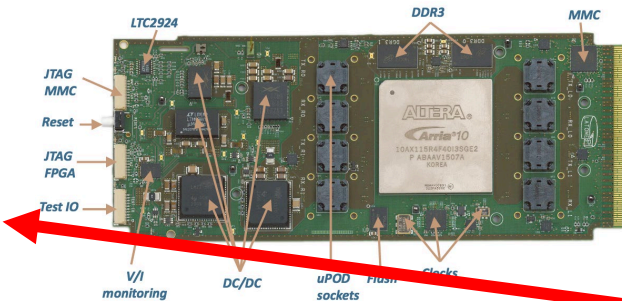
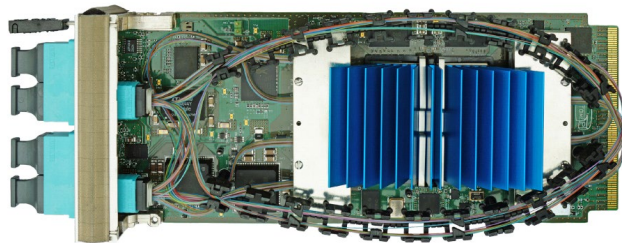
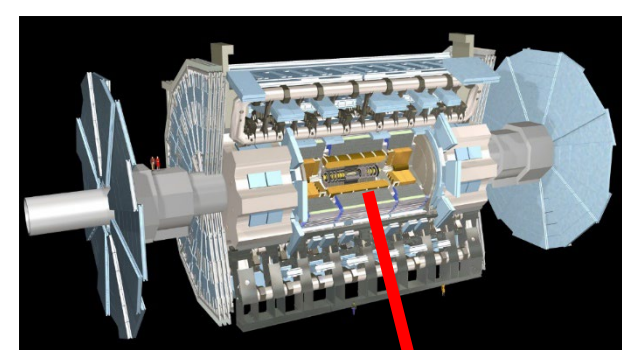
まとめ

- 素粒子実験
 - BigData → 人工知能の応用（情報理工的な理論寄りではなく実践的なツール）
- 既存のツールを**機械学習を取り入れたツール**に置き換える
 - 個人的な感覚的には、5年ぐらい前まではかなり懐疑的（特にATLAS実験グループでは）
 - 今は積極的
 - **圧倒的な改善がときどきある（少々の改善は確実）**
 - 実験データの理解と十分に使えるだけの研究実績・経験
 - 競争相手の存在（ちょっと不甲斐ないですが）
 - 素粒子実験データ
 - グラフネットワークやトランスフォーマは相性がいい
- 素粒子実験のデータ解析用の**基盤モデル（Foundation Model）**の可能性



素粒子屋さんとしては **「素粒子標準模型」の牙城を崩す！**

（もう一枚、宣伝）



カロメータ関連の研究