

Graph Neural Network pour la reconstruction de traces dans le détecteur ATLAS

Sylvain Caillou

Laboratoire des 2 Infinis - Toulouse (L2IT)

sylvain.caillou@l2it.in2p3.fr

13^e Journées Informatique IN2P3/IRFU, Nov 15 – 17, 2021

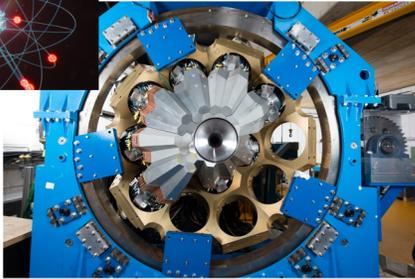


Quelques mots sur moi

- Sylvain Caillou, Ingénieur de Recherche en informatique au CNRS
 - Rentré au CNRS en 2005 en tant que IE
 - A rejoint le L2IT (IN2P3) en Septembre 2020
 - En poste depuis 2008 au LIMSI (INS2I, Institut des Sciences de l'Information et de leurs Interactions) à Orsay
 - Spécialiste des algorithmes d'intelligence artificielle et des bibliothèques associées
 - Formateur en *machine learning* à l'IFSeM (formation CNRS Île de France), membre du comité d'expert pour l'établissement du nouveau catalogue de formation informatique de l'IFSeM 2022-2026

Le L2IT - Toulouse, dernier né des laboratoires IN2P3

Physique nucléaire



GANIL

Ondes gravitationnelles



VIRGO

données

données

L2IT

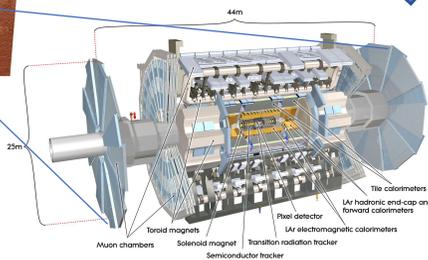
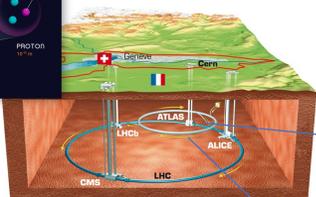
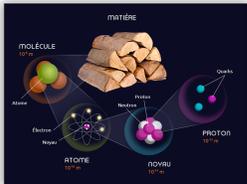
Ressources de calcul



Retour d'expériences / beta tests



Physique des particules



ATLAS

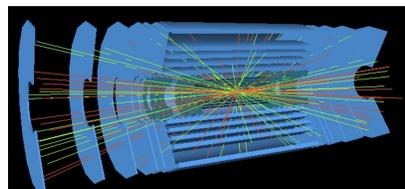
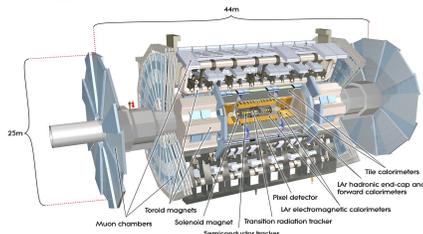
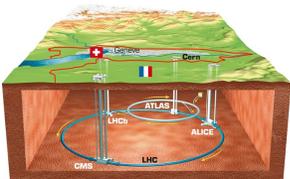
données

Laboratoire des 2 Infinis – Toulouse
Activité principales:

- ⇒ Analyse de données en physique
- ⇒ Développement d'algorithmes innovants, basés entre autres sur le machine learning

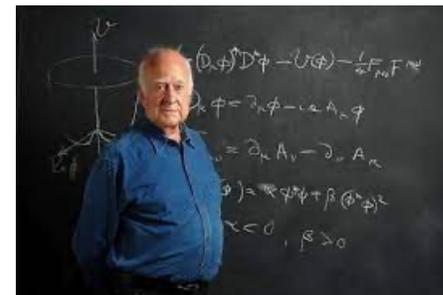
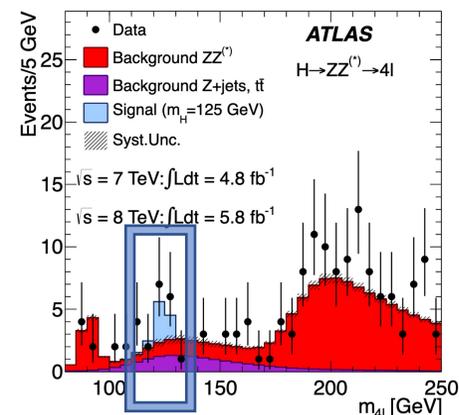
Le challenge de computing de HL-LHC

Large Hadron Collider (LHC) au CERN



Collision de paquet de proton à très haute énergie pour créer des particules
40 millions de croisement de paquets de protons par seconde !

Découverte du boson de Higgs en 2012



La prochaine grande étape pour la physique du Higgs au LHC

Etudier les propriétés du boson de Higgs!

HL-LHC => Intensité du faisceau multiplié par un facteur 10,
démarré en 2027

=> Augmentation du volume et de la *complexité* des données

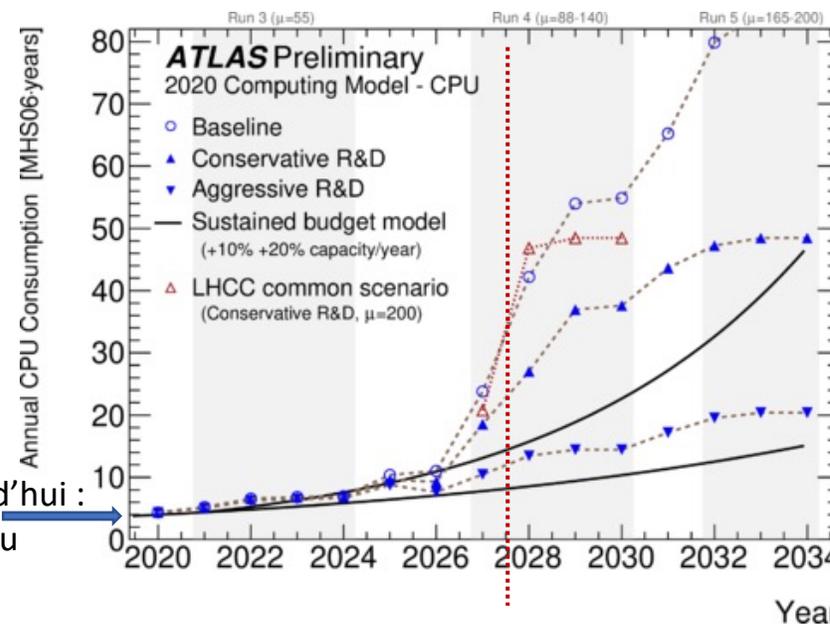
Avec notre modèle de *computing* actuel, il manquera un facteur important dans les ressources de calcul et de stockage

=> **nécessité de changements importants**

=> **au risque de limiter le programme de physique**

Ressources utilisées aujourd'hui :

- O(10⁶) de cœurs en continu
- O(10¹²) octet de stockage



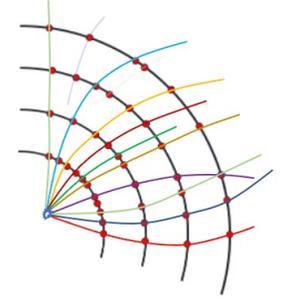
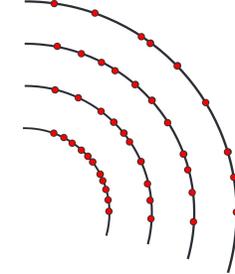
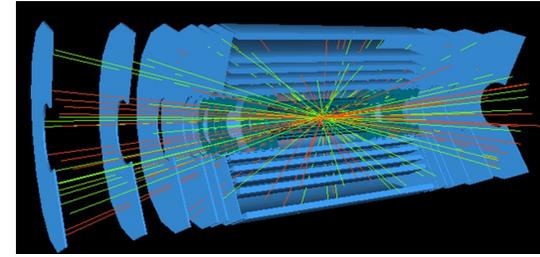
Démarrage de HL-LHC

Le défi de HL-LHC pour la reconstruction de traces dans ATLAS



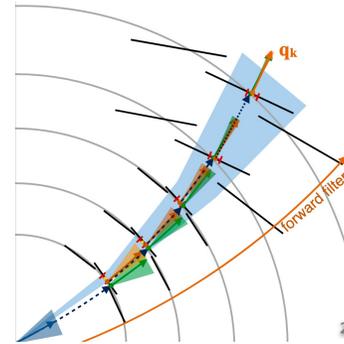
Les particules laissent des « Hits » dans les modules du détecteur

- Trouver les trajectoires (traces) des particules chargées produites dans les collisions p-p
- Nécessaire pour étudier les propriétés des particules produites et faire de la physique



Reconstruction actuellement opérée par des algorithmes basés sur des filtres de Kalman

- Estimation des paramètres de la trajectoire hélicoïdal
- Très bonne performance et optimisé depuis des années
- La partie la plus coûteuse en complexité de calcul (donc en ressources CPU) dans la reconstruction d'un événement



1. propagate p_{k-1} and its covariance C_{k-1} :

$$q_{k|k-1} = f_{k|k-1}(q_{k-1|k-1})$$
$$C_{k|k-1} = F_{k|k-1} C_{k-1|k-1} F_{k|k-1}^T + Q_k$$

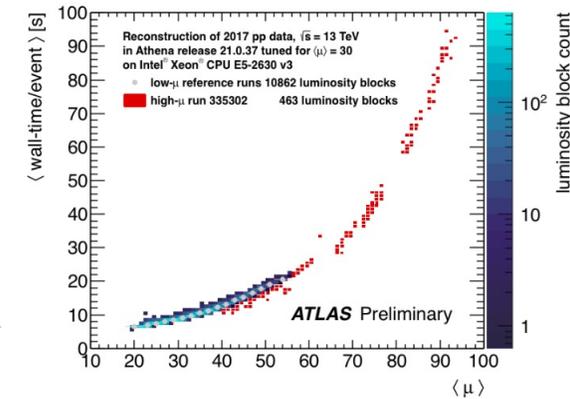
with $Q_k \sim$ noise term (M.S.)

2. update prediction to get $q_{k|k}$ and $C_{k|k}$:

$$q_{k|k} = q_{k|k-1} + K_k [m_k - h_k(q_{k|k-1})]$$
$$C_{k|k} = (I - K_k H_k) C_{k|k-1}$$

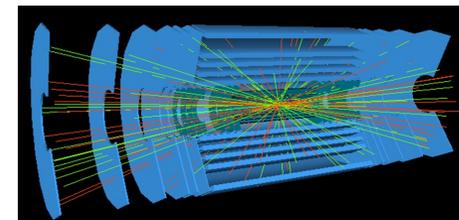
with $K_k \sim$ gain matrix :

$$K_k = C_{k|k-1} H_k^T (G_k + H_k C_{k|k-1} H_k^T)^{-1}$$

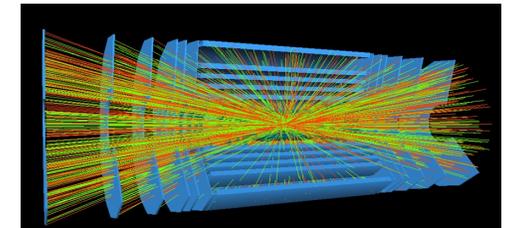


Démarrage de HL-LHC en 2027

- La combinatoire va exploser (pileup $\sim 20 \Rightarrow$ pileup ~ 200)
 - Augmentation très importante du volume et de la complexité des données (10^5 hits et 10^4 traces par évènement)
 - Les algorithmes actuels ne suffiront pas
- \Rightarrow Il nous faut un algorithme fondamentalement nouveau**



LHC (pileup ~ 20)

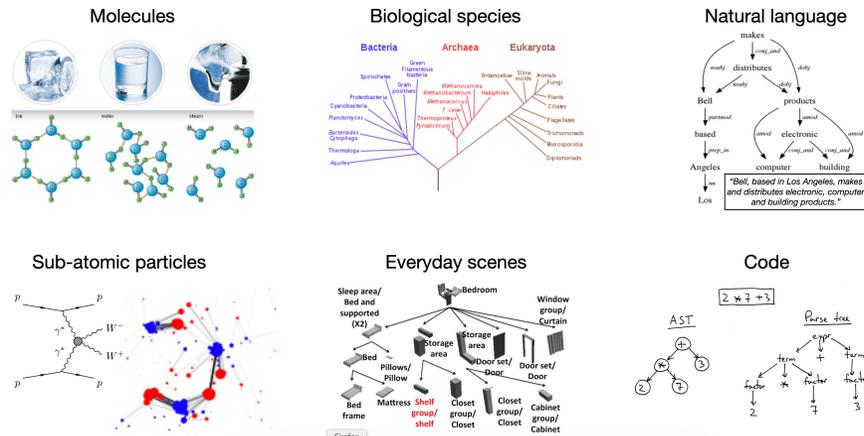


HL-LHC (pileup ~ 200)

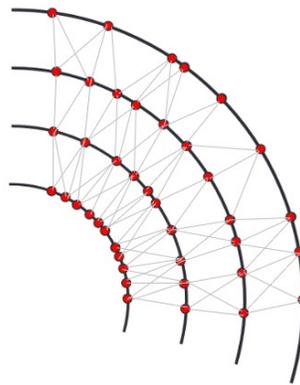
Imaginer un nouvel algorithme basé sur l'IA pour répondre au défi de HL-LHC

Peut-on utiliser l'IA pour aider à reconstruire les traces ?

Beaucoup de systèmes complexes sont structurés et peuvent être représentés sous forme de graphe

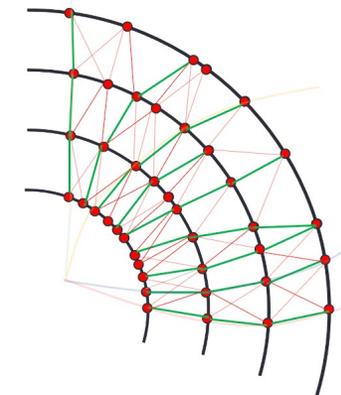


Représentation des données (les hits) sous forme de graphe



Nœuds du graphe = Hits du tracker
Arcs du graphe = Connexions considérées comme possibles

Ce que l'on voudrait faire :
Classifier les arcs du graphe



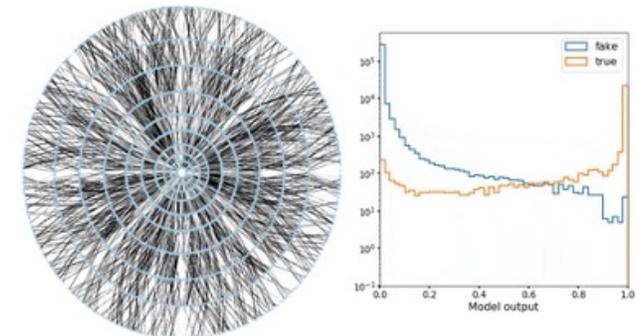
- **Score haut = forte probabilité** que l'arc appartient à une trace
- **Score faible = faible probabilité** que l'arc appartienne à une trace

Exa.Tkrx project (2019 -)

« **Graph Neural Networks for Particle Reconstruction in High Energy Physics detectors** » [arXiv:2003.11603](https://arxiv.org/abs/2003.11603)

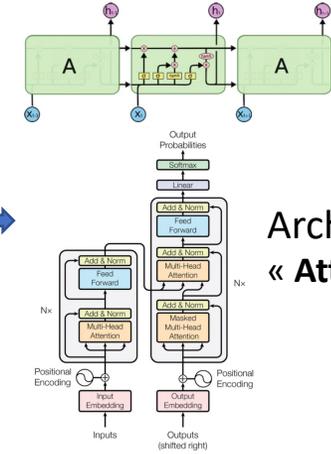
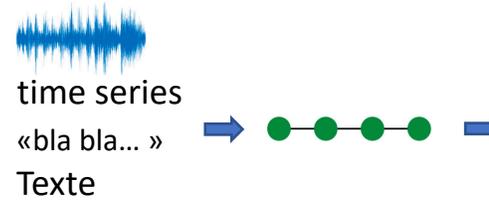
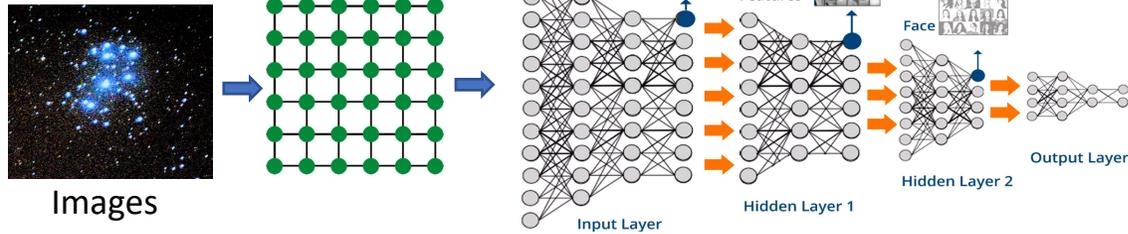
=> Entraîner des réseaux de neurones avec une architecture adaptée les « **Graph Neural Networks** » (GNNs)

=> Preuve de principe sur le dataset du TrackML Data Challenge

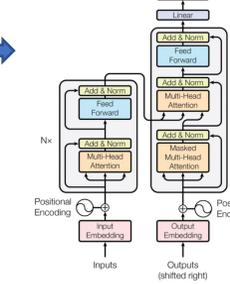


Les Graph Neural Networks (GNNs)

Convolutional Network CNN (~2012)



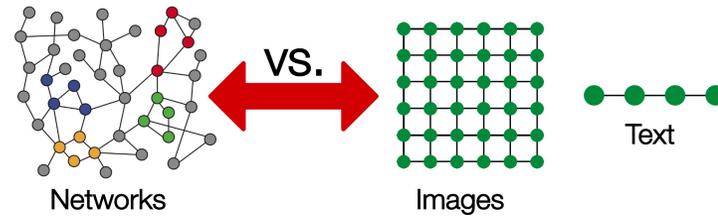
Architectures Transformer « Attention Is All You Need » (2017)



- Graphes permettent de représenter des données éparses avec des relations
- Les GNNs apprennent à partir des features des éléments du graphe mais également à partir des relations qui relient ces éléments entre eux

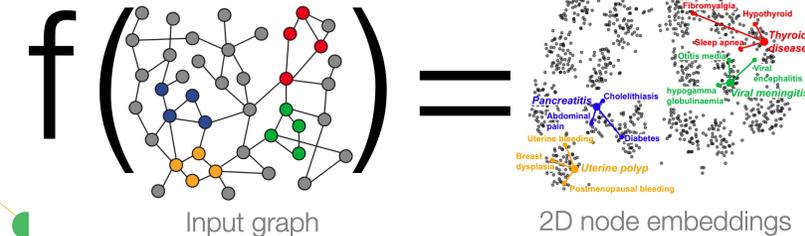
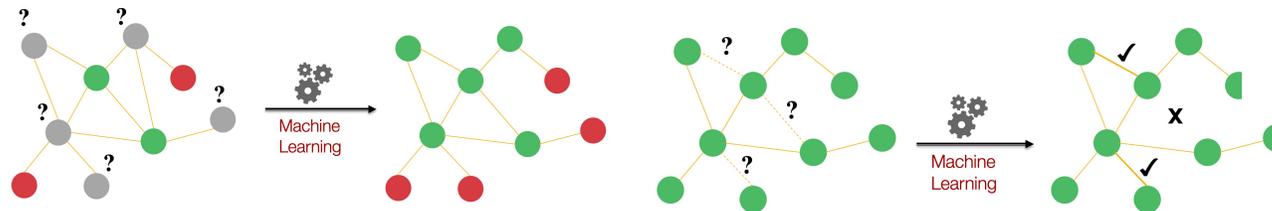
Apprentissage sur les graphes difficile

- Taille arbitraire et structure topologique complexe
- Pas de notion d'ordre ou de point référent



Quelles tâches avec un GNN ?

- Classification de nœuds
- Classification des arcs
- Prédiction de liens
- Classification de graphes



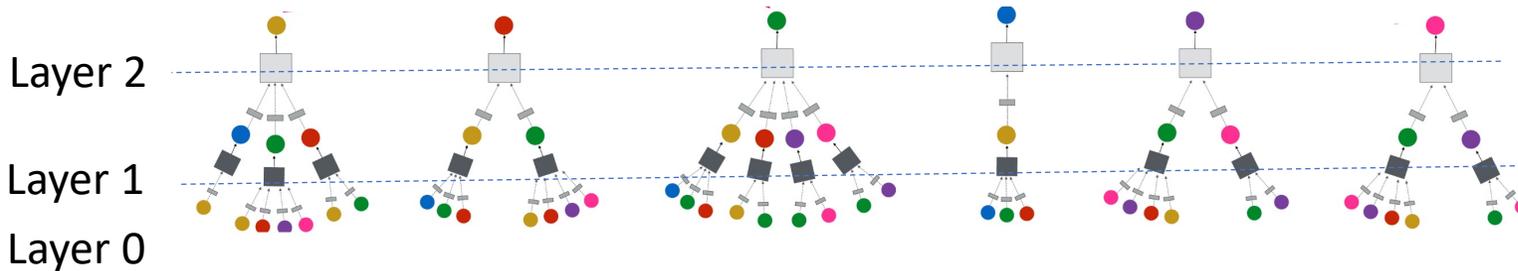
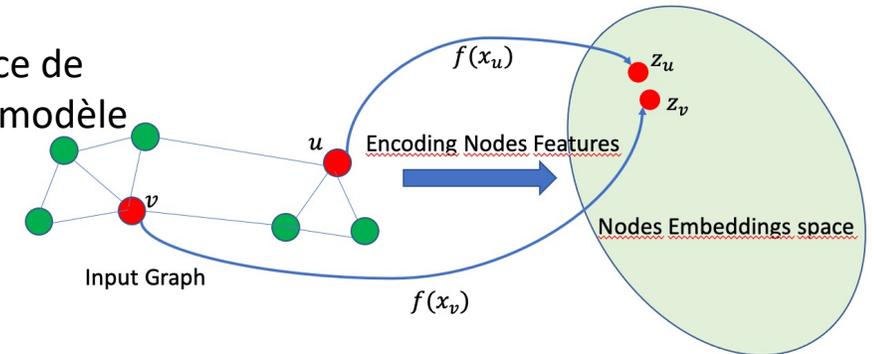
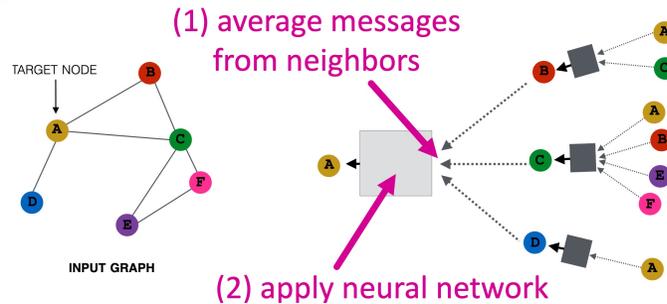
Les Graph Neural Networks (GNNs)

Encoder les nodes (ou les edges) du graphe de manière à ce qu'ils soient proches dans l'espace de projection si ils sont de la même classe et linéairement séparables par la dernière couche du modèle

Pour i allant de 1 à N message passing:

Pour chaque nœud et/ou pour chaque arc du graphe:

- 1) Agréger les informations du voisinage
- 2) Appliquer un réseaux de neurones multicouches non linéaire



Trainable weight matrices (i.e., what we learn)

$$h_v^{(0)} = x_v$$

$$h_v^{(l+1)} = \sigma \left(W_l \sum_{u \in N(v)} \frac{h_u^{(l)}}{|N(v)|} + B_l h_v^{(l)} \right), \forall l \in \{0, \dots, L-1\}$$

$$z_v = h_v^{(L)}$$

Final node embedding

Des bibliothèques python pour construire des architectures GNN existent :

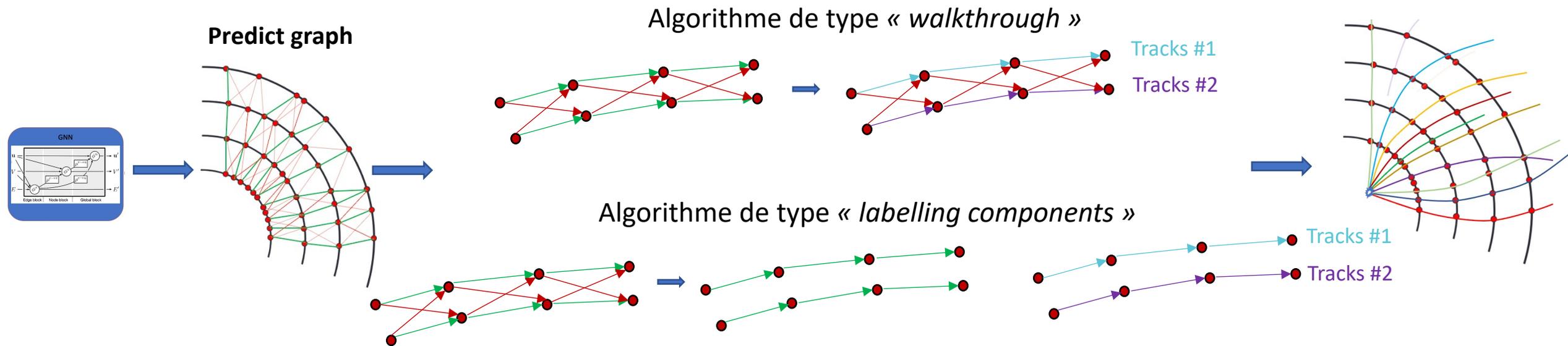
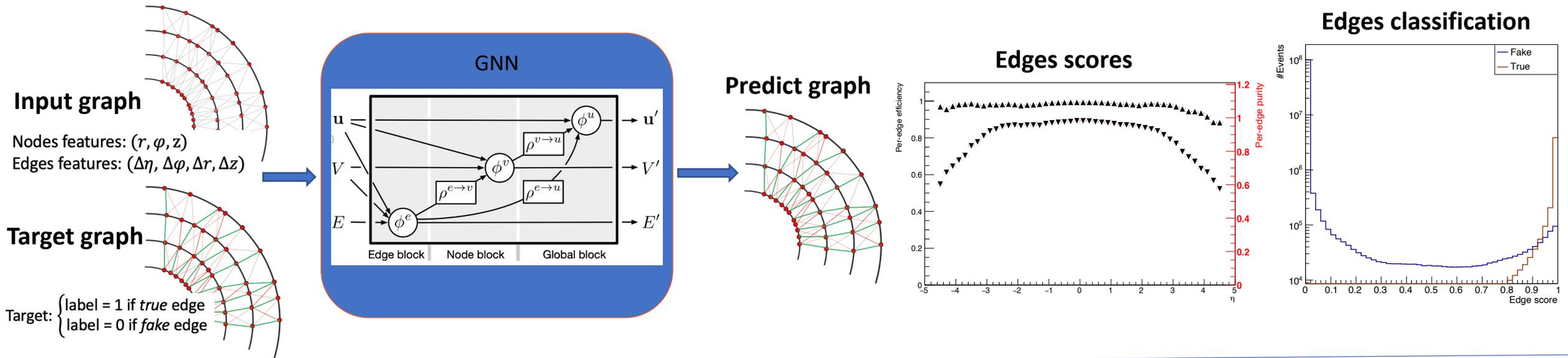
GraphNet (Tensorflow)

Deep Graph Library (DGL) (Tensorflow, Pytorch)

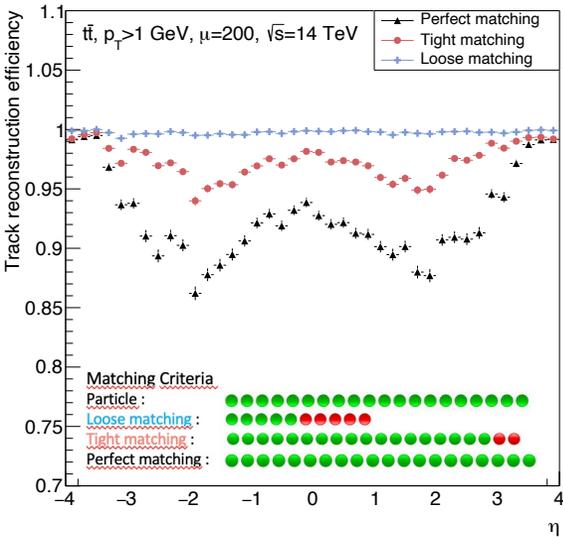
Pytorch Geometrics (PyG) (Pytorch)



Vers un algorithme de reconstruction basé sur les GNNs



Premiers résultats et nouveaux challenges



=> Premier résultats présentés en session plénière à vCHEP 2021



« Towards a realistic track reconstruction algorithm based on graph neural networks for the HL-LHC » (2021)

[Catherine Biscarat](#), [Sylvain Caillou](#), [Charline Rougier*](#), [Jan Stark](#), [Jad Zahreddine](#)
[arXiv:2103.00916](#) [physics.ins-det]

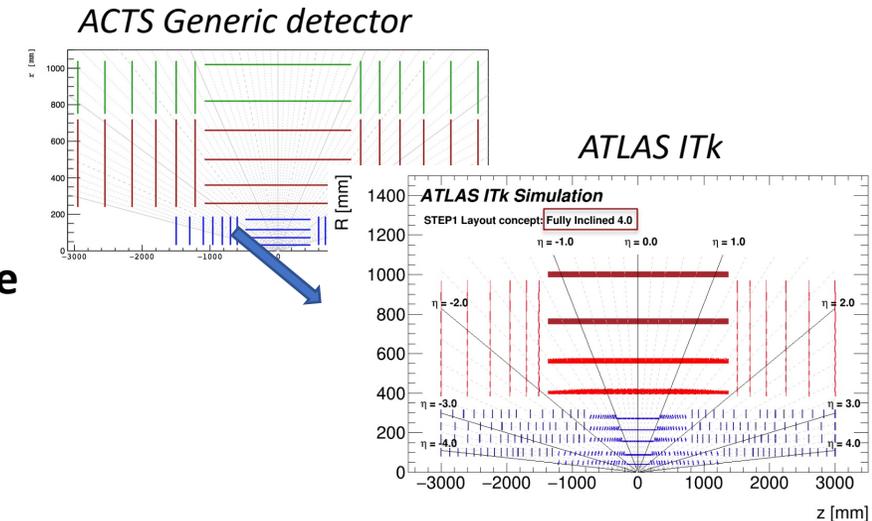
Depuis on est passé à des simulations de données hautement réalistes :

- pileup 200 + diffusion multiple + Ionisation
- Détecteur ITk (futur inner tracker d'ATLAS)

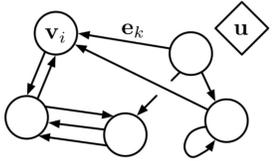
=> Graphes avec $O(10^6)$ edges avec une haute densité et une topologie complexe

Le problème est dur à résoudre

=> On a besoin d'améliorer nos modèles

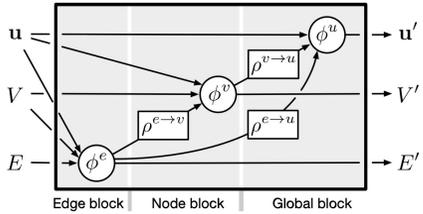
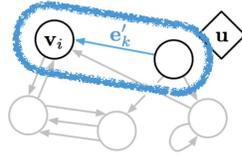


Comprendre et tester différentes architectures GNN



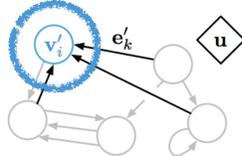
Edge block

For each edge, e_k, v_{s_k}, v_{r_k}, u , are passed to an "edge-wise function":
 $e'_k \leftarrow \phi^e(e_k, v_{r_k}, v_{s_k}, u)$



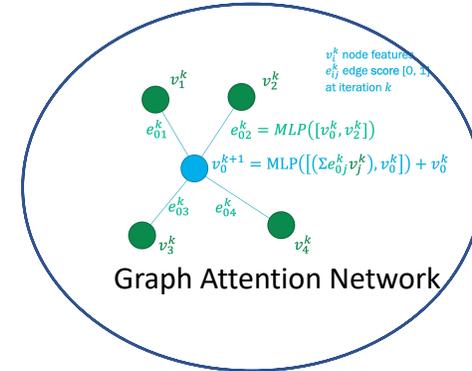
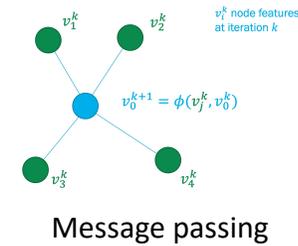
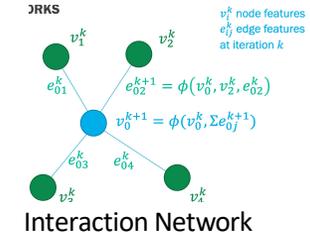
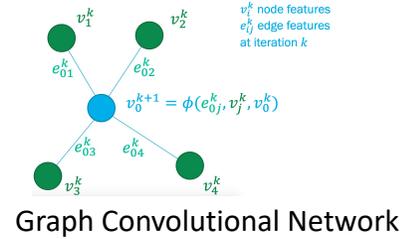
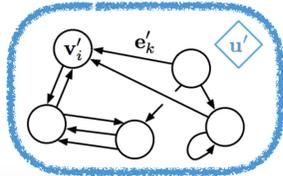
Node block

For each node, e'_i, v_i, u , are passed to a "node-wise function":
 $v'_i \leftarrow \phi^v(e'_i, v_i, u)$



Global block

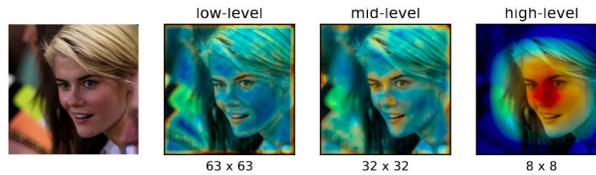
Across the graph, \bar{e}', \bar{v}', u , are passed to a "global function":
 $u' \leftarrow \phi^u(\bar{e}', \bar{v}', u)$



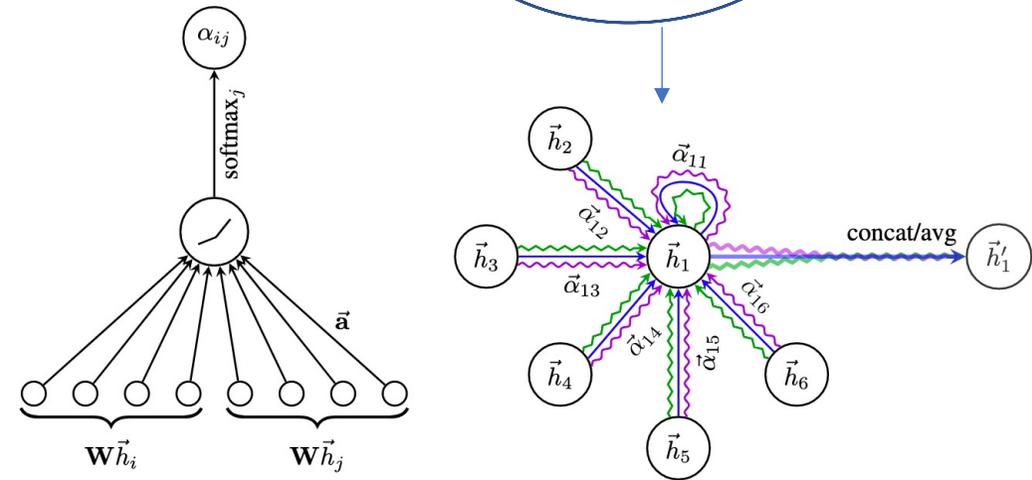
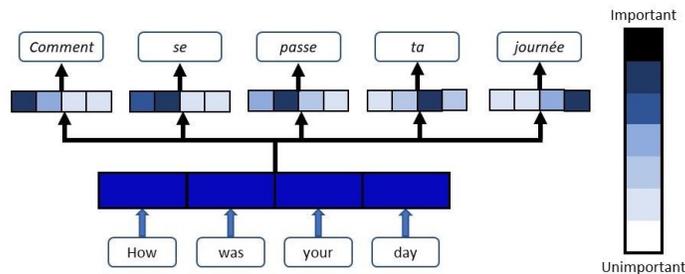
Beaucoup plus encore ...

Modèles bioinspirés => Mécanismes d'attention

Attention spatiale dans les architectures CNN pour la classification ou la segmentation d'images



Attention séquentielle ou temporelle dans les architectures RNN ou les modèles transformer pour le traitement du langage naturel ou les times series



Graph Attention Network (GAT)

[Veličković et al., arXiv:1710.10903 \[stat.ML\]](https://arxiv.org/abs/1710.10903)

Out of Memory Error

Récurrance du GNN

Besoin d'accéder au N message passing dernier états du graphe pour mettre à jour les arcs => Nécessité de les garder en mémoire

$$N \text{ edges} \left[\begin{array}{c} N \text{ edge features} \\ \end{array} \right] \times N \text{ message passing}$$

Rétropropagation du gradient

Chaque arc définit un graphe computationnel
Nécessité de garder en mémoire tous les états cachés à la sortie des couches linéaires

$$N \text{ edges} \left[\begin{array}{c} N \text{ edge features} \\ \end{array} \right] \times (5 \times N \text{ MLP edge layers} + 4 \times N \text{ MLP edge layers} \times N \text{ message passing})$$

Par exemple avec : message passing = 8, N edge features = 128, N edges = $O(10^6)$, N MLP edge layers = 2
On a : $2 \times ((2 \times 8) \times 10^6 \times 128 \times 8 + (2 + 4 \times 8 + 2 + 1) \times 10^6 \times 128 \times 8) \Rightarrow \mathbf{O(100 \text{ Gb})!}$

Techniques investiguées

- IBM TFLMS et Pytorch LMS (swapping entre CPU host et GPU)
- Automatic Mix Precision (AMP) (passer en précision float16 pour les couches linéaires)
- Checkpointing (recalculer le gradient à la volée au moment de la backpropagation)
- stochastic training
- random sampling

Environnement de développement du projet

Construction des graphes

Boost Graph (C++)



Entraînement GNNs

TensorFlow (Graphnet)

PyTorch (DGL, Pytorch Geometrics)



Opérations sur les graphes

NetworkX



Graph-tool



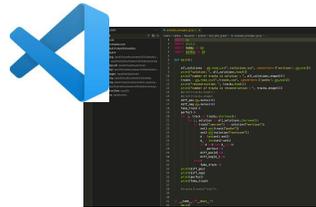
cuGraph

Visualisation entraînement

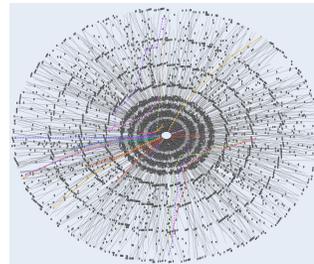
Tensorboard



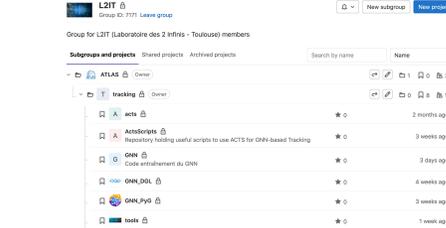
Developpement dans IDE (VSCode)



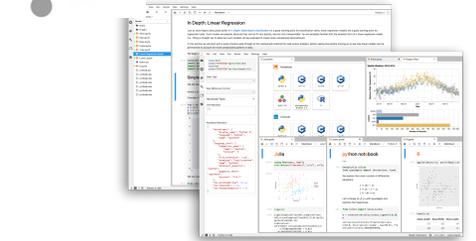
Visualisation données massives



Gestion de version



Analyse de données Prototype de code



Fermes CPU et GPU



Besoins R&D spécifiques L2IT



2xRTX 8000 + Host CPU 1To RAM

4xV100 + Host 1To RAM

MERCI LE CC !

Travaux actuels

- Amélioration des performances des modèles GNN
- Amélioration de l'algorithme de reconstruction (régression sur les trajectoires hélicoïdales des particules)
- Comparaison en terme de temps de calcul et de performances de reconstruction avec le kalman d'Athena (ATLAS)
- Intégration dans ACTS (C++, boost graph, onnx)



MERCI ;-)