

GNNs for the High Granularity Calorimeter at CMS & Online Implementation Ideas

AG THINK II

Matthieu Melennec¹

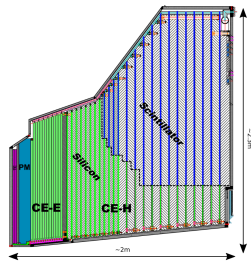
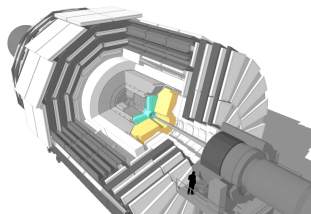
¹Laboratoire Leprince-Ringuet, Ecole Polytechnique, Institut Polytechnique de Paris, CNRS Nucléaire & Particules, Palaiseau, France

16th of October, 2025



High-Luminosity LHC (HL-LHC):

- ▶ More rare events (Higgs production and BSM physics)
- ▶ Increased reconstruction complexity (up to 200 Pile-up events)
- ▶ CMS High Granularity Calorimeter (HGICAL):
 - ▶ New CMS end-cap sampling calorimeter
 - ▶ High granularity: 6M channels on 47 layers
 - ▶ Si and Scintillator based

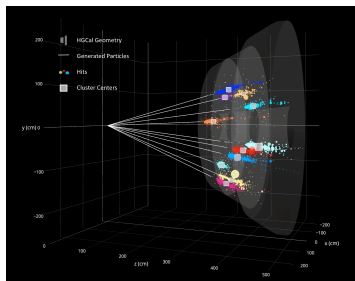


Point cloud data:

- ▶ Points $P_i \in \mathbb{R}^{k \geq 3}$: Euclidean coordinates + $(k - 3)$ “colours”
- ▶ **Unordered, sparse** and with **variable size**

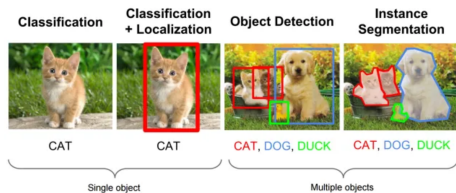
HGCAL output: Point clouds

- ▶ Hits: 3D points with energy measurement and timing
- ▶ Variable granularity
- ▶ Graph convolution promising approach



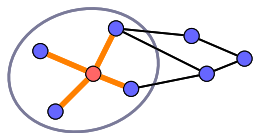
Convolutional Neural Networks:

- ▶ Excellent at classification and segmentation tasks
- ▶ Identifies geometric patterns



How to generalise the success of CNNs to point-cloud data?

- ▶ Graph convolution



Formalism:

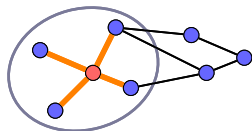
Aggregator: symmetric and normalised (e.g. mean/max) combines all messages

$$x_v^{(t+1)} = \gamma_{\theta_\gamma} \left(x_v^{(t)}, \square_{w \in \mathcal{N}(v)} \phi_{\theta_\phi} \left(x_v^{(t)}, x_w^{(t)}, e_{vw} \right) \right)$$

Update function:
combine messages with
own features

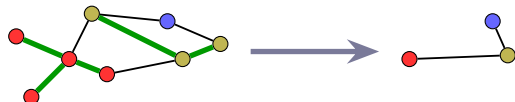
Message function:
collects neighbour features

Gilmer et al., *Neural message passing for quantum chemistry*, 2017



Aim: Coarsen graph to increase the range of the convolution

1. **Selection** (or clustering): Select which nodes to pool (e.g. by selecting edges to “collapse”)
2. **Reduction**: Combine features of pooled nodes (using max or sum pooling)
3. **Connection**: update adjacencies (inherited or dynamic)



Grattarola et al. *Understanding Pooling in Graph Neural Networks*, 2024

Optimal Graph Convolution for particle IDentification Efficient algorithms for event reconstruction in particle detectors

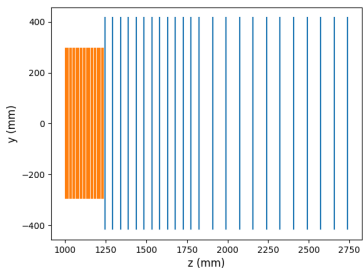
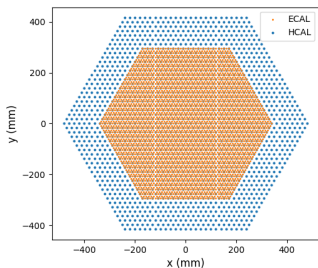
- ▶ Reduction of graph construction complexity
- ▶ Segmented implementation
- ▶ Optimising the network design and adapt it to the electronic implementation (FPGA...)
- ▶ Multi-task (Online and Offline CMS HGAL reconstruction, Hyper-Kamiokande DSNB discrimination)



Funded by the Agence Nationale de la Recherche (ANR), ANR-21-CE31-0030

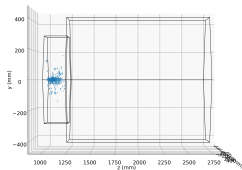
Simulate HGCAL-like calorimeter using GEANT4

- ▶ $\sim 10^5$ Si sensors
- ▶ 26 ECAL layers with Pb absorbers
- ▶ 24 HCAL layers with stainless steel absorbers

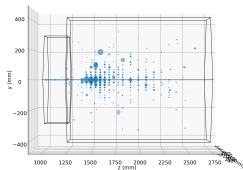


Simulated e^-/γ , π^+ and μ^- events in the detector

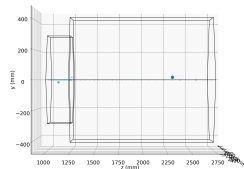
- ▶ Energies 10 GeV to 100 GeV
- ▶ Each hit corresponds to the energy deposited in the detector in the corresponding sensor



e^-/γ event

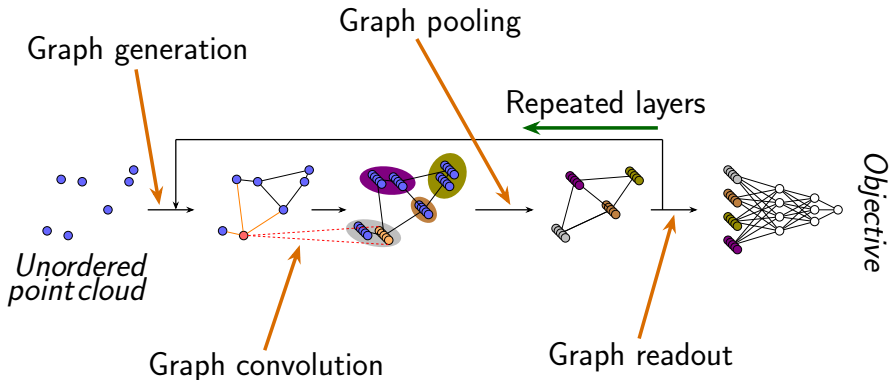


π^+ event

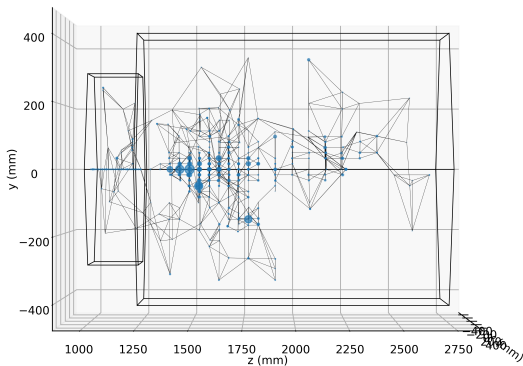


μ^- event

GCNN Structure




- ▶ Build arbitrary edges between sparse, multi-dimensional data-points
- ▶ Typically: k nearest neighbours (KNN)
 - ▶ Ensures geometric locality
 - ▶ Complexity: worst-case = mean = $\mathcal{O}(n^2)$



Particle detectors: Static and known geometry

Pre-compute proximities of sensors

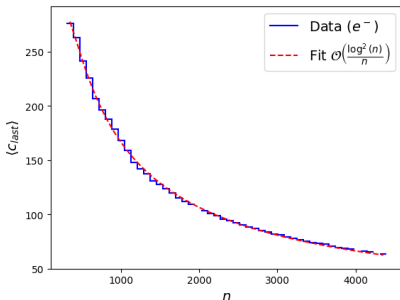
- ▶ For each sensor, order its neighbours by increasing distance in a “proximity table” (PT)

Sensor IDs	Increasing order wrt. metric			
				
1	17	38	...	42
2	75	16	...	68
...				
99	3	98	...	22

- ▶ Arbitrary choice of metric used for ordering (e.g. Euclidean, adding a radially term, correlation...), but no correlation on model performance in our study: take Euclidean distance

PT-KNN: iterate over rows until k neighbours found

- ▶ Worst-case complexity: $\mathcal{O}(n^2)$
- ▶ Reduces best case complexity to $\mathcal{O}(kn)$
- ▶ On average, $\langle c_{\text{last}} \rangle \sim \frac{\log^2(n)}{n}$

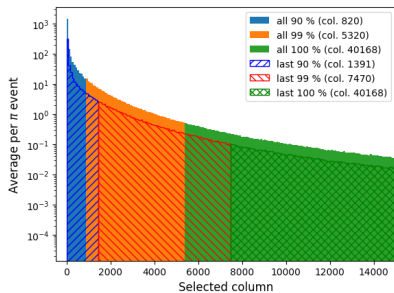
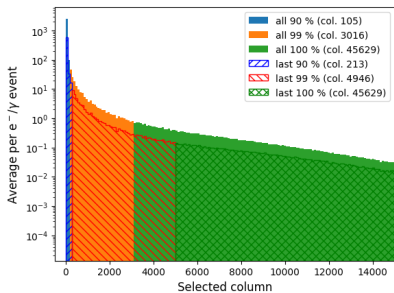


PTs reduce the mean complexity of KNN (in the relevant range) from

$$n^2 \quad \text{to} \quad \log^2(n)$$

Proximity Tables: $10^5 \times 10^5$ entries

- ▶ Can cut PT to remove rarely explored columns
- ▶ Allows FPGA implementation



We obtain graphs:

▶ **Nodes** v :

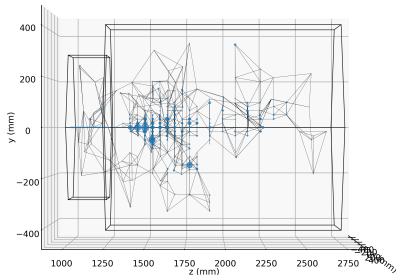
- ▶ Sensor energy x_v
- ▶ Position \vec{u}_v

▶ **Graph-level features** (pid, energy...)

Radial symmetry in detector \Rightarrow
Positions \vec{u}_v carried as “hidden
features”, not used in convolution

▶ **Edges** e_{vw} :

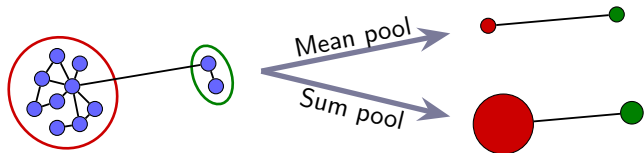
- ▶ End nodes v, w
- ▶ Length $d(v, w) = \|\vec{u}_v - \vec{u}_w\|$



$$x_v^{(t+1)} = \square_{w \in \tilde{\mathcal{N}}(v)} \text{Leaky-ReLU} \left(\Phi_\theta \left[x_v^{(t)} \ x_w^{(t)} \ d(v, w) \right] \right)$$

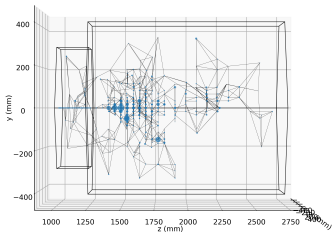
- ▶ Message function Φ_θ : Linear combination with trainable weights θ
 $\Phi_\theta \in \mathbb{R}^{2n \times (2n+1)}$, i.e. doubles number of features
- ▶ Aggregator \square : Feature-wise pooling (classification: max, regression: mean)
- ▶ Update function γ : Self-loop (i.e. aggregate with message from itself)

1. **Selection with Treclus:** Collapse all edges shorter than a threshold ε
 - ▶ Choice of ε using the number of resulting nodes: Convolution doubles n° features \Rightarrow pooling halves n° nodes
2. **Reduction:** Combine nodes v in cluster \mathcal{C}
 - ▶ Feature-wise pool $\{x_v\}_{v \in \mathcal{C}}$ (classification: max, regression: sum)
 - ▶ Choose at random a destination node in $\{\vec{u}_v\}_{v \in \mathcal{C}}$

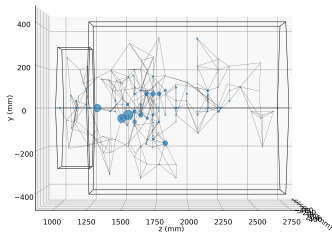


3. **Connection:** Inherited adjacency from nodes $v \in \mathcal{C}, w \in \mathcal{C}'$ neighbours $\Rightarrow \mathcal{C}, \mathcal{C}'$ neighbours

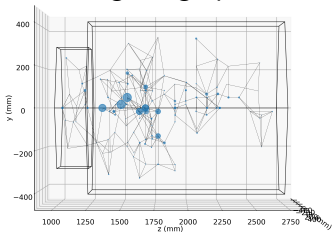
Example Pooling



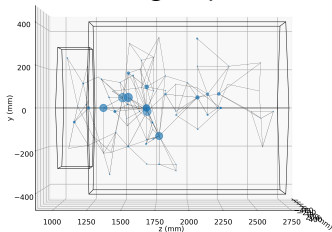
Original graph



Pooling step 1



Pooling step 2

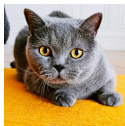


Pooling step 3

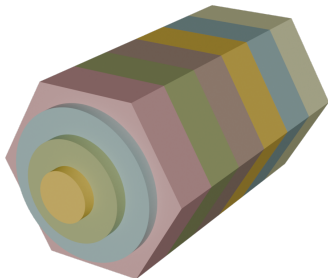
Readout problematic:

- ▶ Need to flatten graph structure as input for an MLP
- ▶ Can be tricky to keep graph structural information
 - ▶ No order for nodes
 - ▶ No order for edges
- ▶ Need a consistent approach

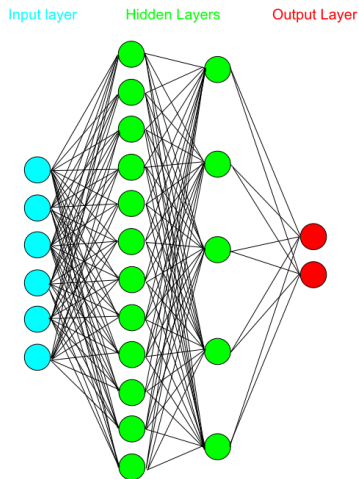
Random order of readout unintelligible →



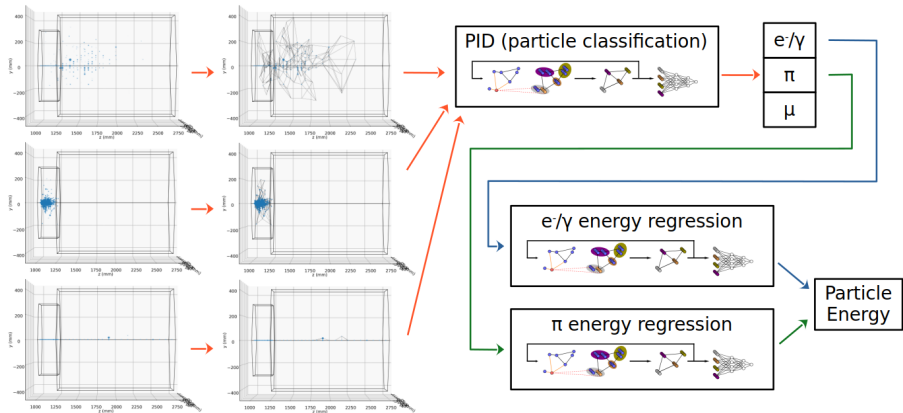
- ▶ Known geometry: embed graph back into its geometry
- ▶ Detector sliced up in readout regions that respect rotational symmetry
- ▶ Pool features within the same region (max or sum)
- ▶ Flatten in consistent order



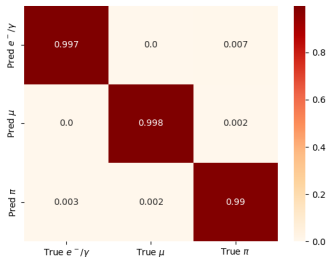
- ▶ Fully connected MLP
- ▶ 5-6 hidden layers
- ▶ Leaky ReLU activation
- ▶ Output size: 3 (PID) or 1 (Energy regression)



Full Pipeline

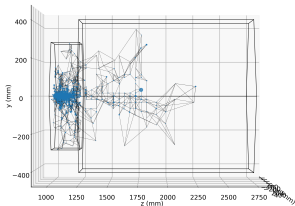


Pipelines have 3 CP layers, 6 hidden MLP layers $\sim 10^4$ parameters
Readout granularity adapted to task

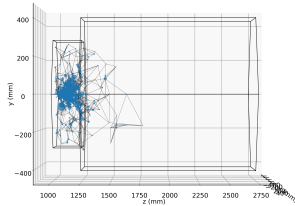


- ▶ Classify $e^-/\gamma, \mu, \pi$ with $E \in [10, 100]$ GeV
- ▶ Balanced set of 10^5 events
- ▶ State of the art performance
- ▶ Some difficult PID tasks

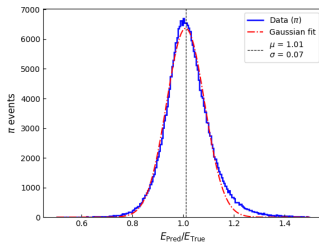
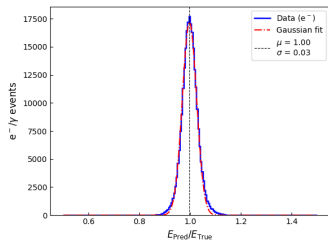
e^- induced Hadronic jet



Early showering π



- ▶ Trained on 2×10^6 graphs (75% training)
- ▶ Regression performance conform to detector
- ▶ e^-/γ better precision than π : different sampling fractions and physics
- ▶ Asymmetry of tails: detector properties



Energy resolution given by:

Stochastic fluctuations
in shower development

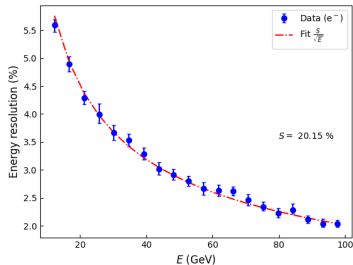
Noise from read-
out electronics

$$\left(\frac{\sigma\left(\frac{E_{\text{Pred}}}{E_{\text{True}}}\right)}{\left\langle \frac{E_{\text{Pred}}}{E_{\text{True}}} \right\rangle} \right)^2 = \frac{S}{\sqrt{E_{\text{True}}}} + \frac{N}{E_{\text{True}}} + C$$

Systematic noise
(e.g. dark noise...)

► Noise not emulated:

$$\frac{\sigma}{\mu} \propto \frac{1}{\sqrt{E}}$$



- ▶ Graph convolution powerful tool for HEP data
- ▶ Recover state of the art results
- ▶ Algorithmic optimisation allows online implementation (e.g. FPGAs)

Perspectives:

- ▶ More difficult PIDs
- ▶ Bigger energy range
- ▶ Pile-up Segmentation
- ▶ Extension to other detectors (e.g. diffuse supernovae background in Hyper-Kamiokande)

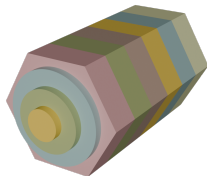
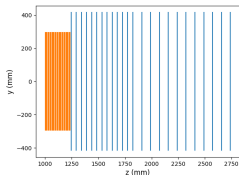
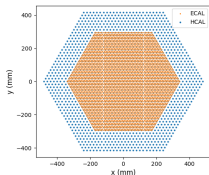
Part II

Online Implementation

Detector Geometry Table

Save the detector geometry in a static table ($M \times 6$)

Cell ID	x (mm)	y (mm)	z (mm)	Layer	RO sector
1	0	0	1000.32	1	0
2	10	0	1000.32	1	0
...
66807	-60	-190.526	1168.98	19	42
...
M	240	-415.692	2740.3	50	106

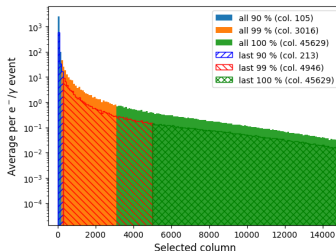


Proximity Table

As above, ranks each other sensor by increasing distance.
Static, size ($M \times M - 1$)

Cell ID	Increasing distance \rightarrow			
1	61	11	...	130609
2	62	13	...	121463
...			...	
M	136069	136019	...	3570

Can be reduced to $\sim 10^3$



Data is saved in a $(N_{\max} \times 2)$ table

Cell ID	Energy (MeV)
183	0.3849
932	2.3671
...	...
N	1.4372
$N + 1$	0
...	...
N_{\max}	0

- ▶ Maximum number of nodes N_{\max} .
- ▶ N_{\max} strikes a balance between memory usage and performance.

We represent graphs as $(N_{\max} \times 3 + 3K + F)$ tables

Node ID	Cell ID	Exists	Pooled	Neighbour 1				...	Neighbour K			Node features		
				Cell ID	Exists	Dist (mm)	...		Cell ID	Exists	Dist (mm)	f_1	...	f_F
1	52	True	False	17	True	142.4	...	3528	False	1203.54	0.21	...	1.98	
2	396	False	True	42	True	59.23	...	130	False	501.59	3.67	...	2.3	
...	
N	27	True	True	260	False	10.35	...	602	True	1802.31	0.198	...	0.005	
$N + 1$	0	False	False	0	False	0	...	0	False	0	0	...	0	
...	
N_{\max}	0	False	False	0	False	0	...	0	False	0	0	...	0	

- ▶ Nodes are not ordered (pseudo-random order with a LUT).
- ▶ Maximum of $K > k$ neighbours per node. The minimal value for K is determined by the pooling.
- ▶ F features, corresponding to the features at all layers for skip/residual connections.

The PT-KNN Algorithm

Require: Data table D , Proximity Table P

Ensure: Graph G

- 1: $G \leftarrow$ empty graph
- 2: **for** $X \in D$ **do**
- 3: Add X to G
- 4: **for** $X \in D$ **do**
- 5: **for** $Y \in D$ **do**
- 6: **if** $|\mathcal{N}(X)| == K$ **then**
- 7: Break
- 8: $\mathcal{N}(X) \leftarrow \mathcal{N}(X) \cup \{Y\}$
- 9: **if** $|\mathcal{N}(Y)| < K$ **then** $\mathcal{N}(Y) \leftarrow \mathcal{N}(Y) \cup \{X\}$
- 10: **return** G

Require: Graph G , Convolution layer l

- 1: **for** $X \in G$ **do**
- 2: $f_X^{(l)} \leftarrow 0$
- 3: **for** $Y \in \mathcal{N}(X)$ **do**
- 4: $m \leftarrow \sigma \left(\Theta \left[f_X^{(l-1)}, f_Y^{(l-1)}, d(X, Y) \right] \right)$
- 5: $f_X^{(l)} \leftarrow \square \left(f_X^{(l)}, m \right)$

- ▶ DSP for weights Θ and aggregator \square
- ▶ LUT for non-linearity σ
- ▶ Complexity $N * K$ but can be parallel

Require: Graph G

Ensure: Cluster C

```
1:  $C \leftarrow \{1, \dots, N_{\max}\}$ 
2: for  $X \in G$  do
3:   if  $\neg Pooled(X)$  then
4:     for  $Y \in \mathcal{N}(X)$  do
5:       if  $Exists(Y) \wedge \neg Pooled(Y)$  then
6:          $C[Y] \leftarrow X$  (or pseudo-random choice)
7:          $Pooled(X, Y) \leftarrow \text{False}$ 
8:         Break
9: return  $C$ 
```

$Pooled(X)$ and $Exists(X)$ return the corresponding flags in G

A $Merge(X, C[X])$ function ensure inheritance of $\mathcal{N}(X)$ by $C[X]$

Require: Graph G , Cluster C

```
1: for  $X \in G$  do
2:   if  $C[X] \neq X$  then
3:      $Exists(C[X]) \leftarrow \text{False}$ 
4:     for  $Y \text{ in } \mathcal{N}(C[X])$  do
5:        $\mathcal{N}(Y) \leftarrow (\mathcal{N}(Y) \setminus \{C[X]\}) \cup \{X\}$ 
6:        $\mathcal{N}(X) \leftarrow \mathcal{N}(X) \cup \{Y\}$ 
```

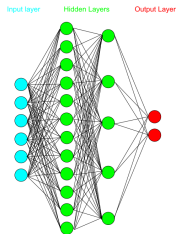
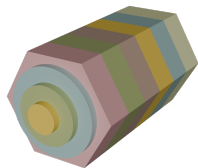
- ▶ Distance from X to Y is either computed with a DSP or inherited from $C[X]$ (i.e. use $d(X, C[X])$)
- ▶ For full graph, complexity K^2

Readout vector has shape $S \times F_{RO}$, with $F_{RO} \leq F$ a subset of all collected features, containing at least the output from the last convolution.

Aggregation within a RO sector done with a DSP

Output from RO is then passed through an MLP (DSP).

Standard implementation (e.g. HLS4ML)



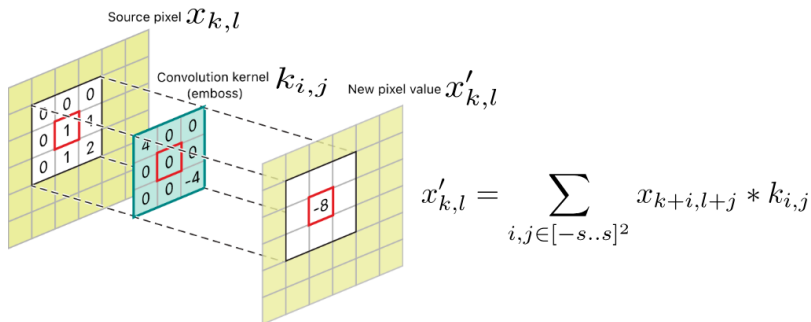
Next steps:

- ▶ We are working on an HLS implementation of this pipeline
- ▶ Evaluate the necessary resources
- ▶ Study the impact of DSP/LUT approximations on performance
- ▶ Study the impact of hyperparameters (K , PT size, distance update in merge ...) on performance

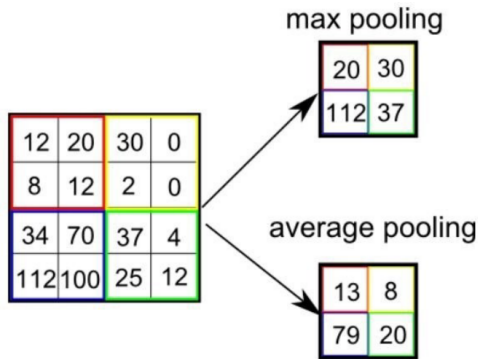
Thank you for listening...
Any questions?

Backups

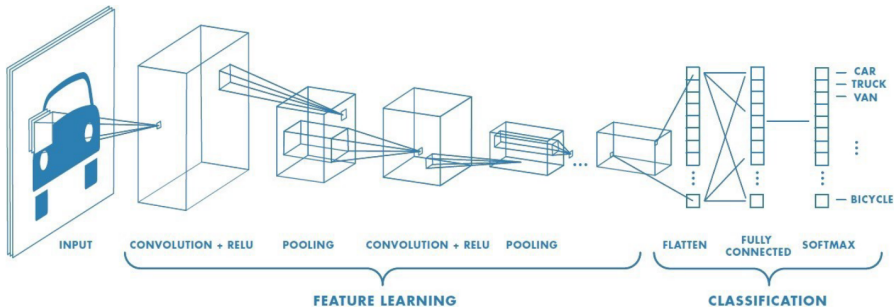
Image Convolution



- ▶ Apply kernel on image (like the convolution filter)
- ▶ Kernel (k_{ij}) is learnable
- ▶ Filter is shared over the whole picture
- ▶ Idea : creating maps of features (one kernel per feature)



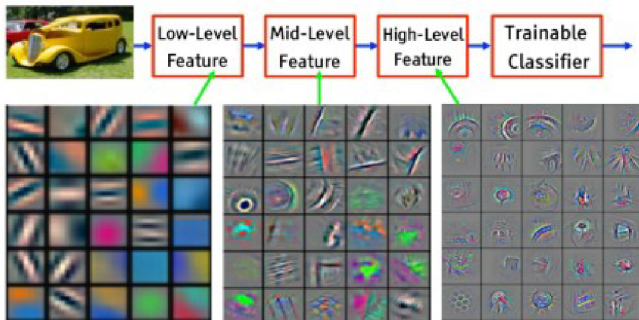
- ▶ Reduce the dimensionality of the feature maps
- ▶ Move to higher level of abstraction
- ▶ Classification: max pool; Regression: mean pool



Network structure :

- ▶ Alternance of convolution & pooling
- ▶ Flattening (sometimes called readout)
- ▶ Multi-layer perceptron

How Does It Work?



- ▶ Feature maps aggregates more and more details to converges to high level recognition patterns
- ▶ Flattened high-level feature map is input for multi-layer perceptron

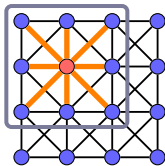
Why Does It Work?

- ▶ The two operations derive naturally from local space:
 - ▶ Euclidean space \Rightarrow Translation invariance; Respected by convolution
 - ▶ Scale-separability \Rightarrow alternated convolution and downsampling
- ▶ Dream complexity
 - ▶ $\mathcal{O}(1)$ parameters par filter (independent of image size)
 - ▶ $\mathcal{O}(n)$ complexity in time per layer (n pixels)

- ▶ Message: $\phi_{\theta_{\phi}}(x_v, x_w, e_{vw}) = x_w * \theta_{\phi_w}$
- ▶ Aggregator: $\square = \Sigma$
- ▶ Every node is self-looped:

$$x_v^{(t+1)} = \sum_{w \in \tilde{\mathcal{N}}(v)} x_w * \theta_{\phi_w}$$

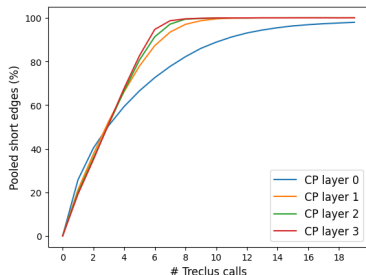
with $\tilde{\mathcal{N}}(v)$ a regular structure containing $\mathcal{N}(v)$ and v



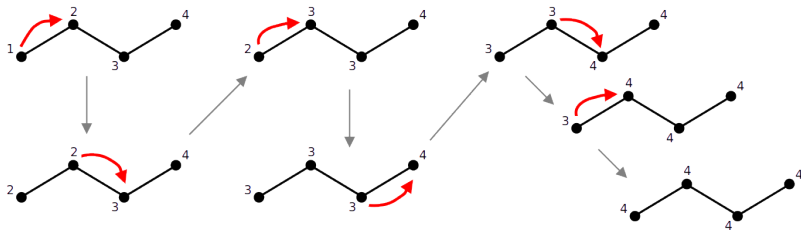
1. Selection with Treclus: Collapse all edges shorter than a threshold ε

- ▶ Choice of ε using the number of resulting nodes: Convolution doubles n° features \Rightarrow pooling halves n° nodes

- ▶ Make a random matching: avoid chain clusters
- ▶ Treclus cannot collapse all short edges: Call multiple times



Chain clusters



Other method of implementing data online:

Cell ID	Energy	Hit
1	1.267	True
2	0	False
...	...	
M	0.344	True

Imposes small changes in graph implementation:

- ▶ G now has size $(M \times 3 + 2K + F)$ and nodes are identified by the address of their cell.
- ▶ Increased memory consumption but simpler node access.