# Pruning and resizing deep neural networks for FPGA implementation in trigger systems at collider experiments

D. Mascione, M. Cristoforetti, A. Di Luca, F. M. Follega, R. Iuppa, A. Saccardo

University of Trento, Fondazione Bruno Kessler, INFN TIFPA

ACAT 2022, Bari 27/10/2022

# *A lot of data*

Collider experiments produce a **huge amount of data**.

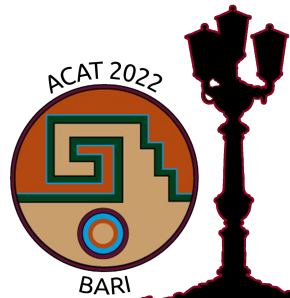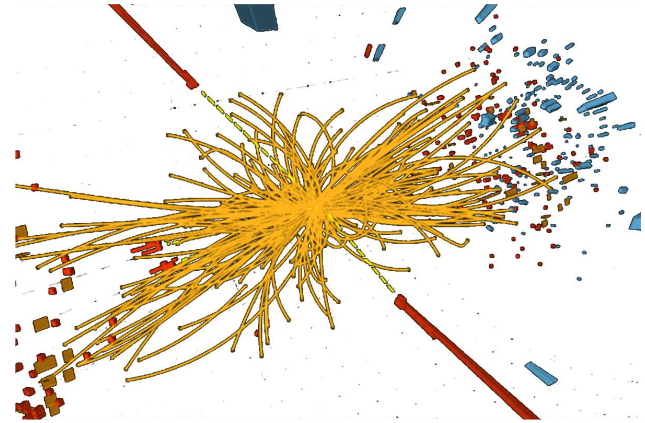At the Large Hadron Collider we have

- one collision every 25 ns ( = **40 Million collisions/sec**)

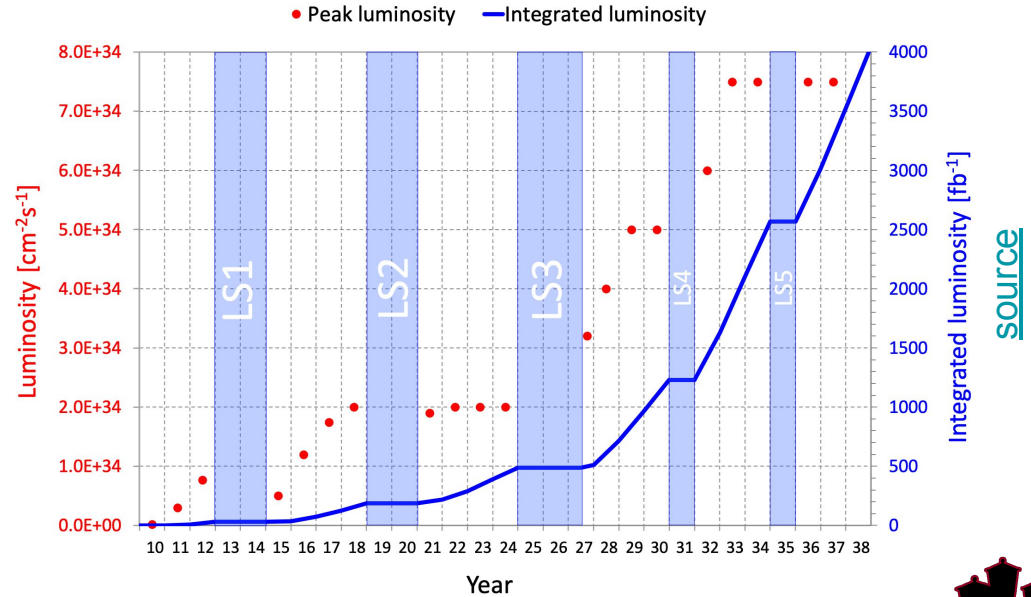- **thousands of particles** emerging from each collision

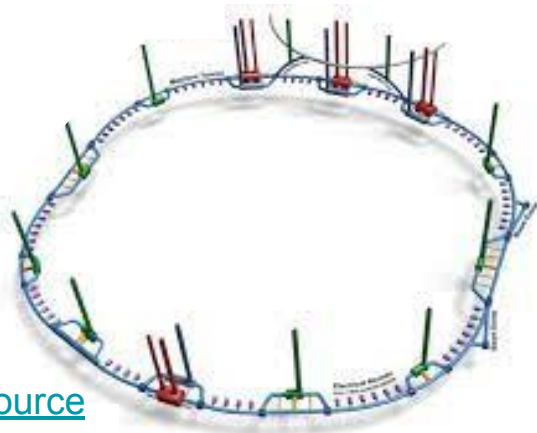- **1 MB of data** recordered at each collision by big detectors

ACAT 2022

BARI

# Increasing data at future colliders

The **HL-LHC** will produce more than 250 fb⁻¹ of data per year and will be capable of collecting **up to 4000 fb⁻¹** (1 fb⁻¹ ~ 100 million million collisions).

At the **FCC-hh** huge amounts of data will be produced (**O(TBytes/s) expected**).

# The trigger system at the LHC

Not all data produced at the LHC are stored: they are first filtered with a trigger chain
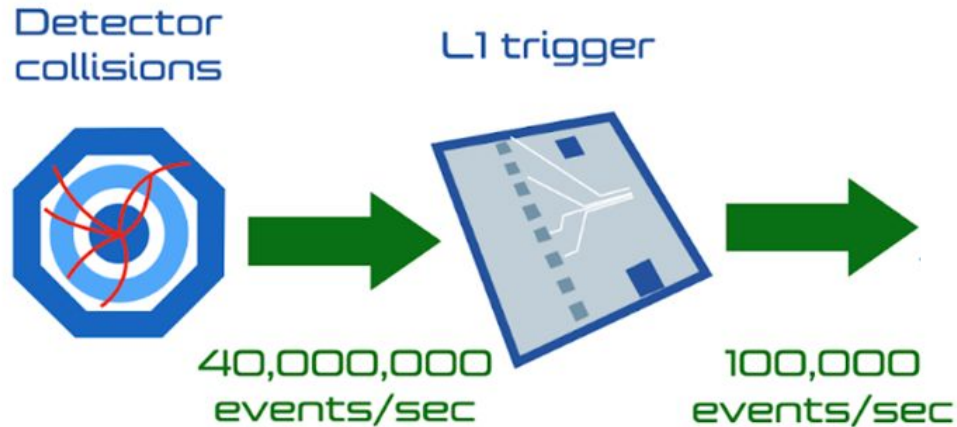
Detector
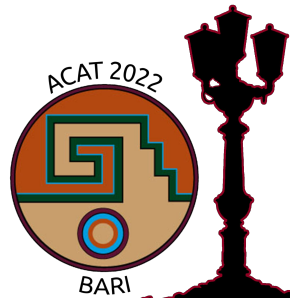collisions

40,000,000
events/sec

source

ACAT 2022

BARI

# The trigger system at the LHC

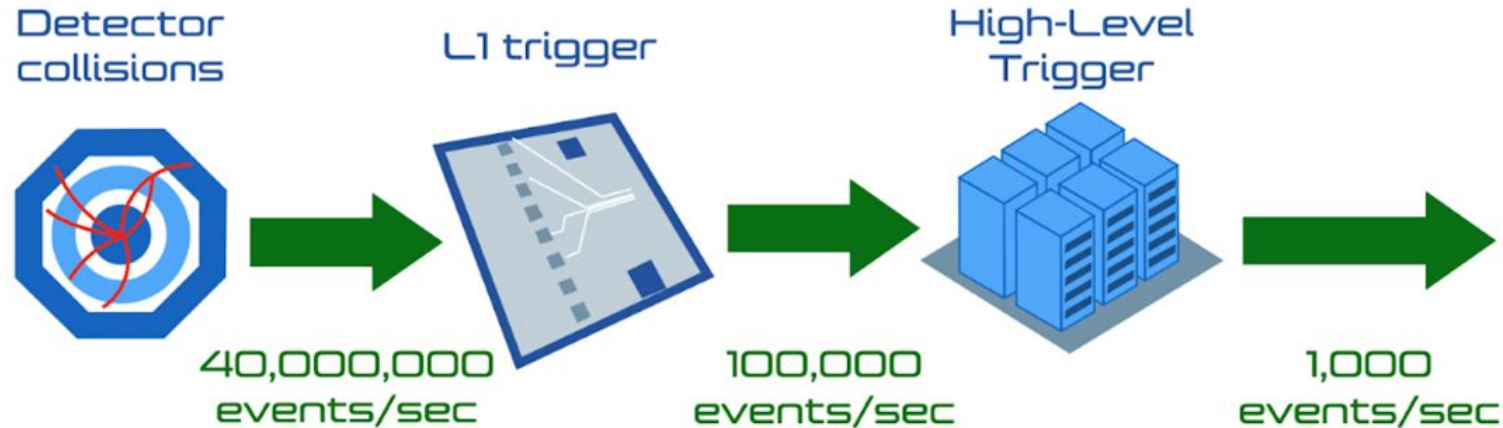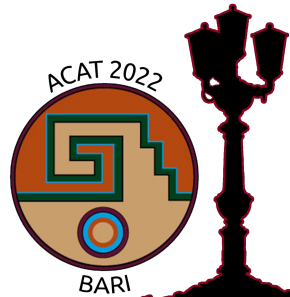Not all data produced at the LHC are stored: they are first filtered with a trigger chain
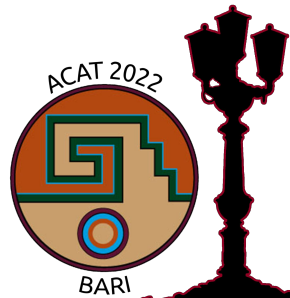


source

# The trigger system at the LHC

Not all data produced at the LHC are stored: they are first filtered with a trigger chain



Detector collisions → L1 trigger → High-Level Trigger →

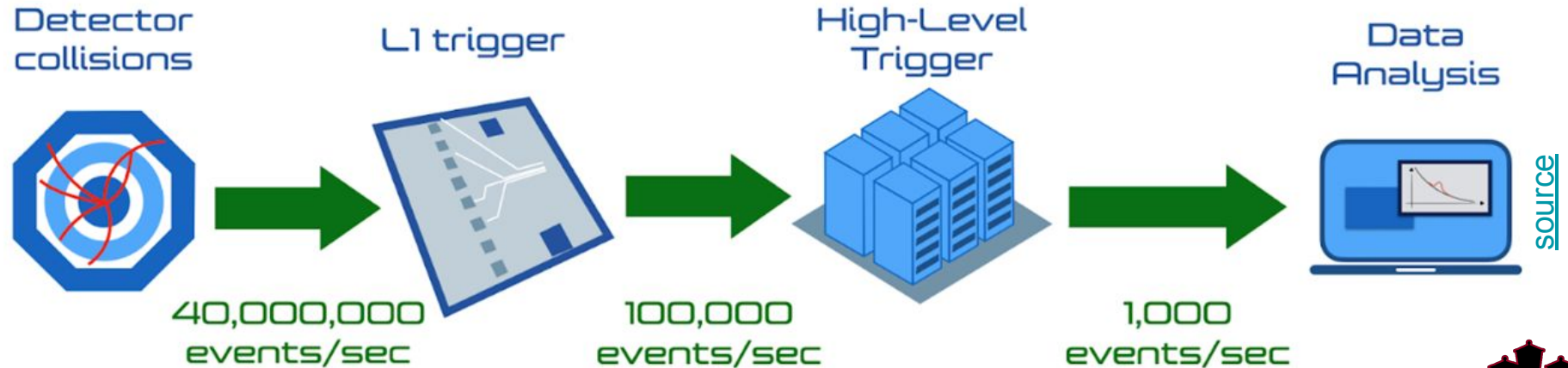40,000,000 events/sec → 100,000 events/sec → 1,000 events/sec
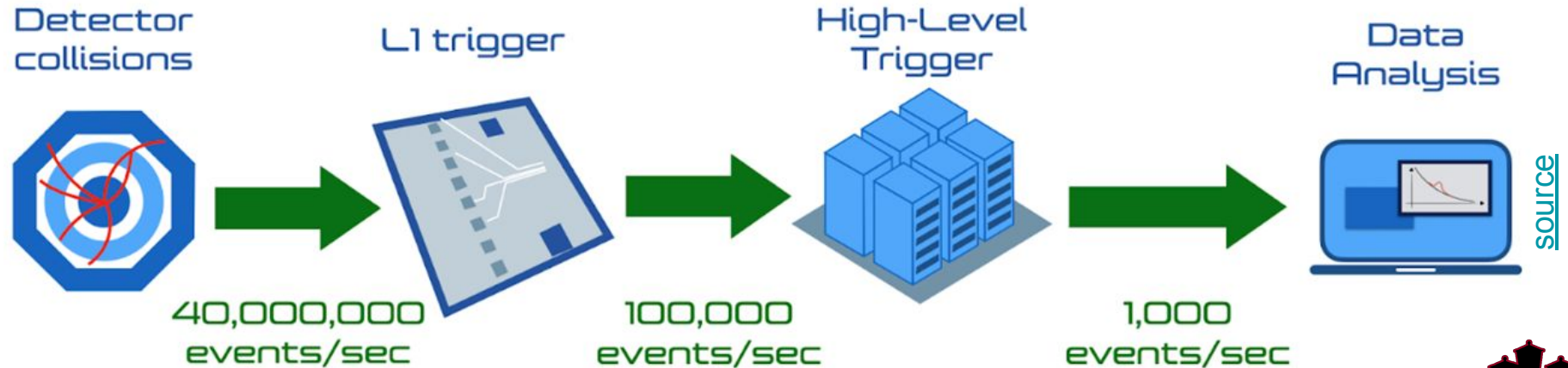
source

ACAT 2022
BARI

# *The trigger system at the LHC*

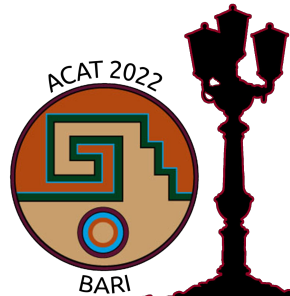Not all data produced at the LHC are stored: they are first filtered with a trigger chain

# The trigger system at the LHC

Not all data produced at the LHC are stored: they are first filtered with a trigger chain
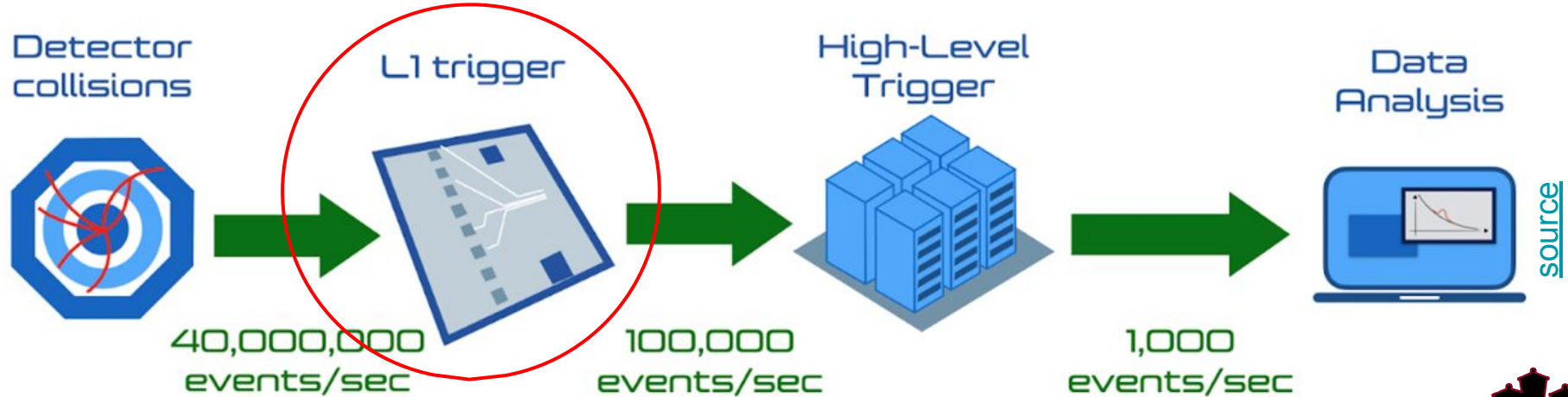


**Detector collisions** → **L1 trigger** → **High-Level Trigger** → **Data Analysis**

40,000,000 events/sec → 100,000 events/sec → 1,000 events/sec

source

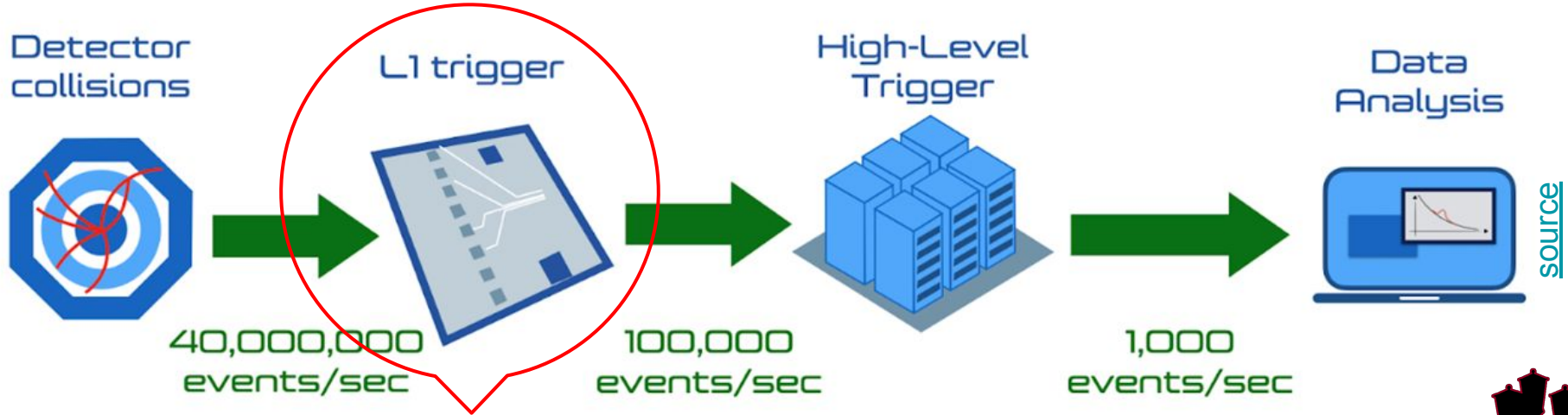⚠️ Events that are discarded by the trigger are **lost**!

# Deep Neural Networks at rescue

Deep Neural Networks can make a **fast event selection** in an extremely dense environment, and can therefore be used where the event selection happens.
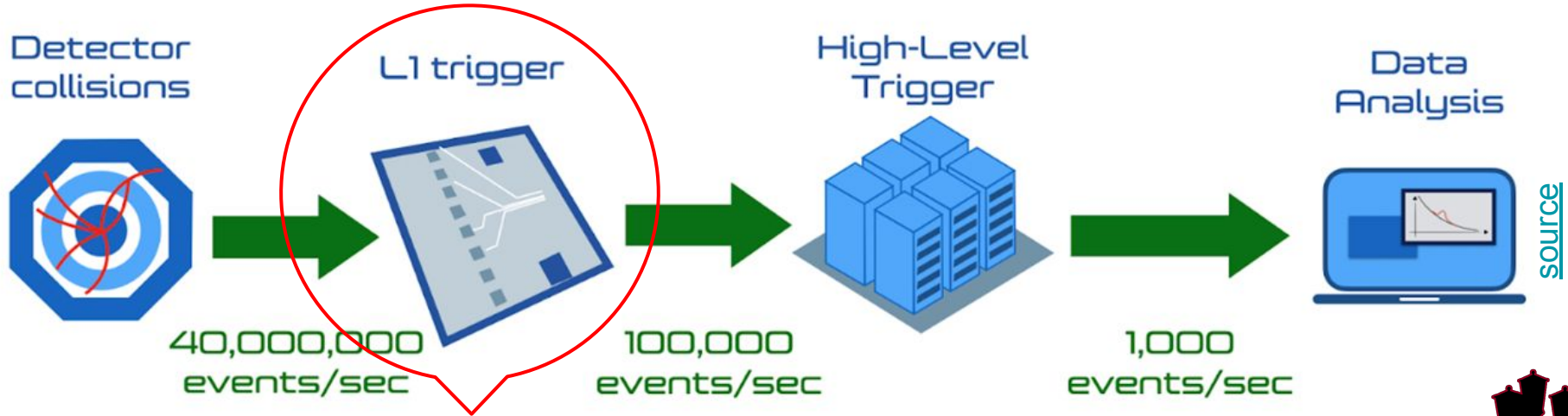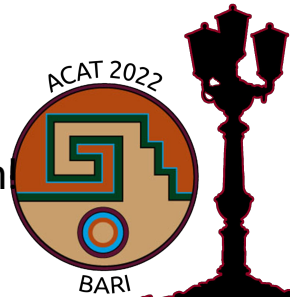


source

# Deep Neural Networks at rescue

Deep Neural Networks can make a **fast event selection** in an extremely dense environment, and can therefore be used where the event selection happens.
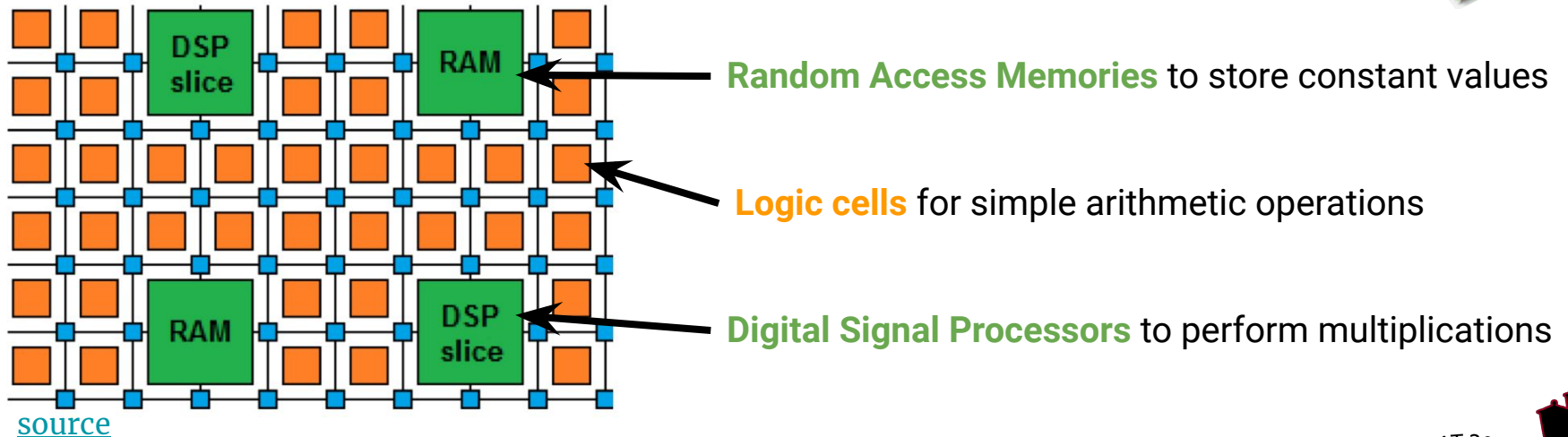


Detector collisions → 40,000,000 events/sec → L1 trigger → 100,000 events/sec → High-Level Trigger → 1,000 events/sec → Data Analysis

source

L1 of data processing typically uses custom hardware with FPGAs

ACAT 2022

BARI

# Deep Neural Networks at rescue

Deep Neural Networks can make a **fast event selection** in an extremely dense environment, and can therefore be used where the event selection happens.



💡 L1 of data processing typically uses custom hardware with FPGAs

💡 Let's run Deep Neural Networks in real-time on FPGAs to improve event selection

# FPGAs

FPGAs (Field-Programmable Gate Arrays) are programmable integrated circuits.



**Random Access Memories** to store constant values

**Logic cells** for simple arithmetic operations

**Digital Signal Processors** to perform multiplications

[source](#)

👉 Depending on the FPGA resources available,
we should know how to **reduce the size** of a network.

# Pruning

One way of **reducing** the size of a neural network is **pruning**.
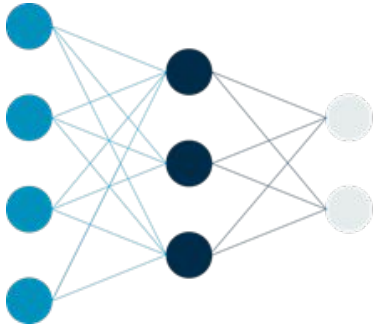
Pruning = **removing** superfluous structure

before pruning

after pruning

source

ACAT 2022

BARI

# The usual pruning scheme

1. Train



Davis Blalock et al., *What is the state of neural network pruning?*, Proceedings of machine learning and systems 2 (2020), pp. 129–146
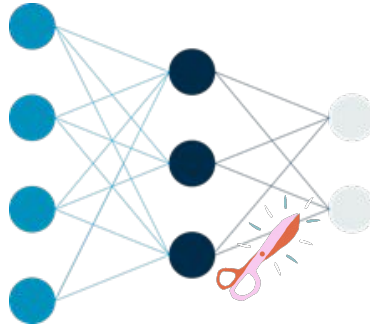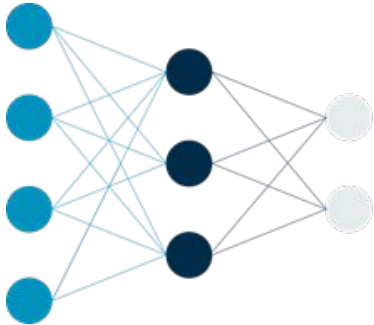
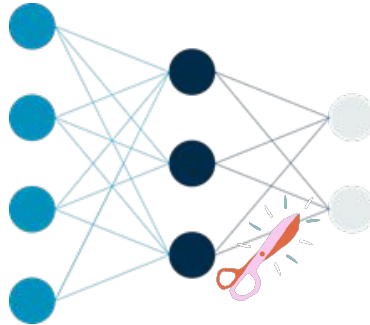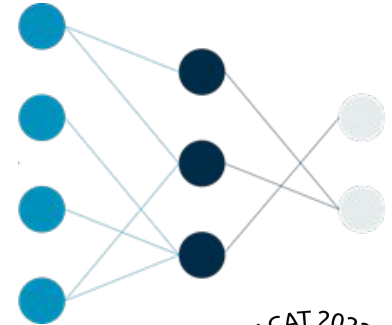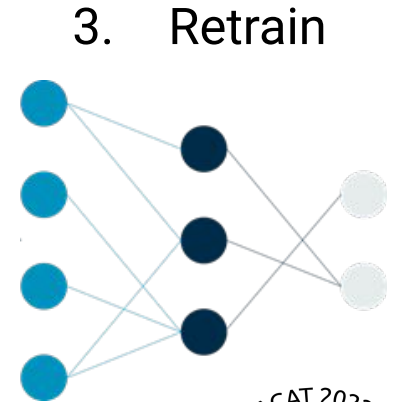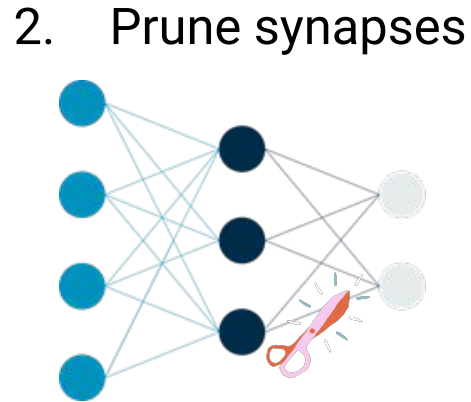# The usual pruning scheme

1. Train

2. Prune synapses

Davis Blalock et al., *What is the state of neural network pruning?*, Proceedings of machine learning and systems 2 (2020), pp. 129–146

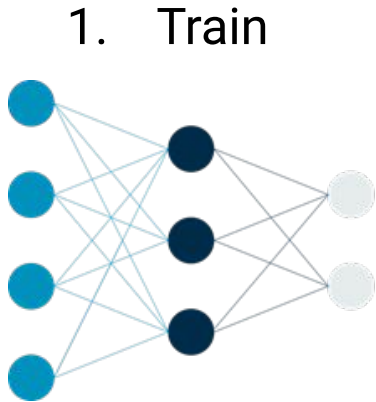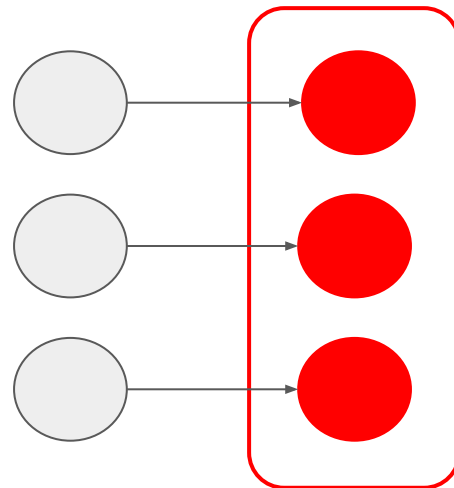# The usual pruning scheme

1. Train    2. Prune synapses    3. Retrain



Davis Blalock et al., *What is the state of neural network pruning?*, Proceedings of machine learning and systems 2 (2020), pp. 129–146

ACAT 2022

BARI

# The usual pruning scheme

Iterate (fine tuning)

1. Train

2. Prune synapses

3. Retrain



Davis Blalock et al., *What is the state of neural network pruning?*, Proceedings of machine learning and systems 2 (2020), pp. 129–146

# AutoPruner: a different pruning strategy

AutoPruner

- it can prune **nodes**

- it prunes **during training**

- the number of nodes to be pruned can be determined by the **user**

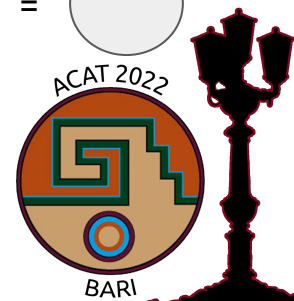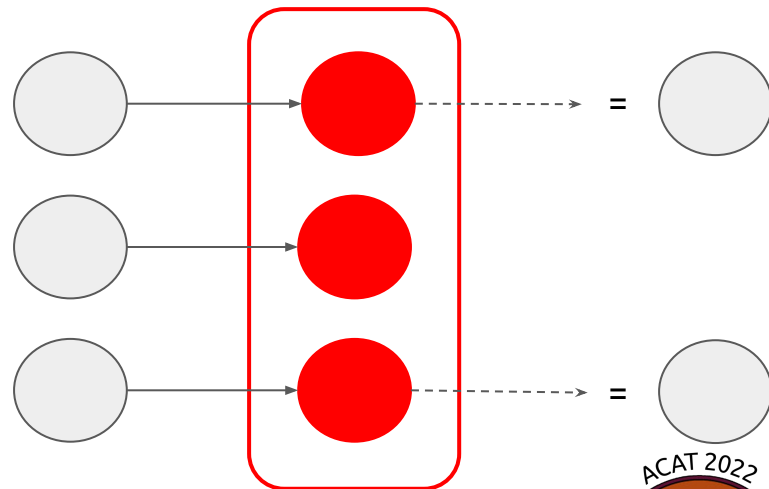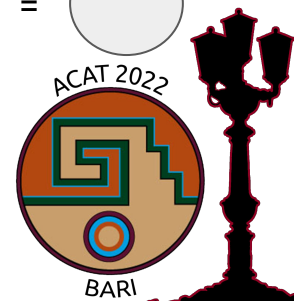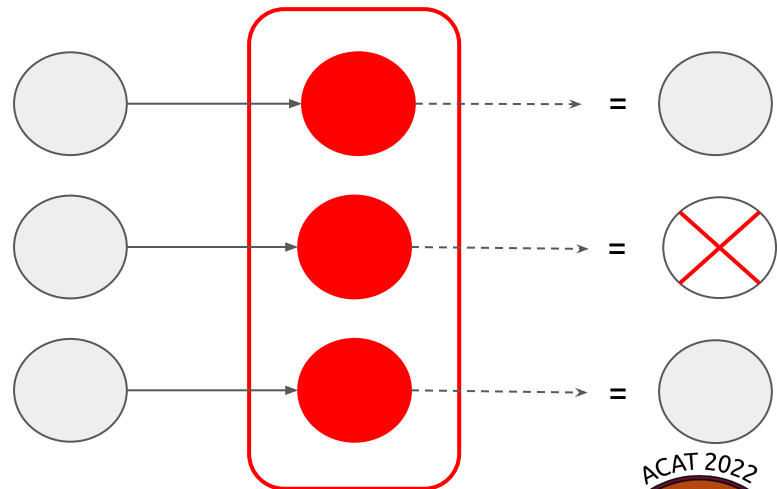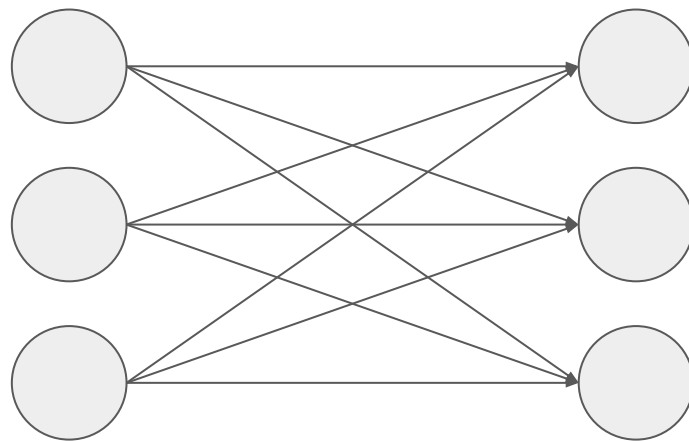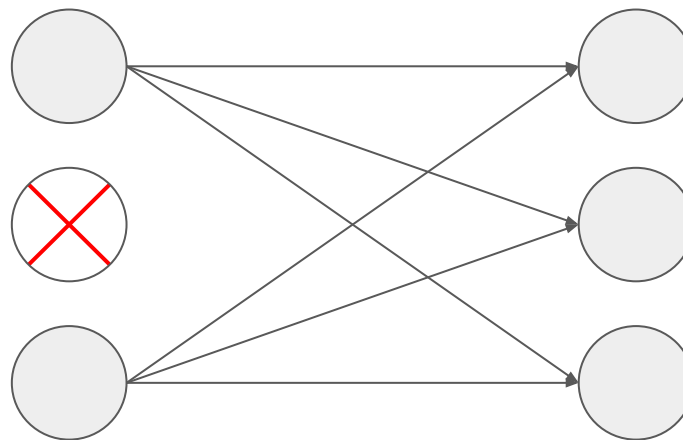- it can determine the most suitable **network architecture**

# AutoPruner: a different pruning strategy

- it can prune **nodes**

- it prunes **during training**

- the number of nodes to be pruned can be determined by the **user**

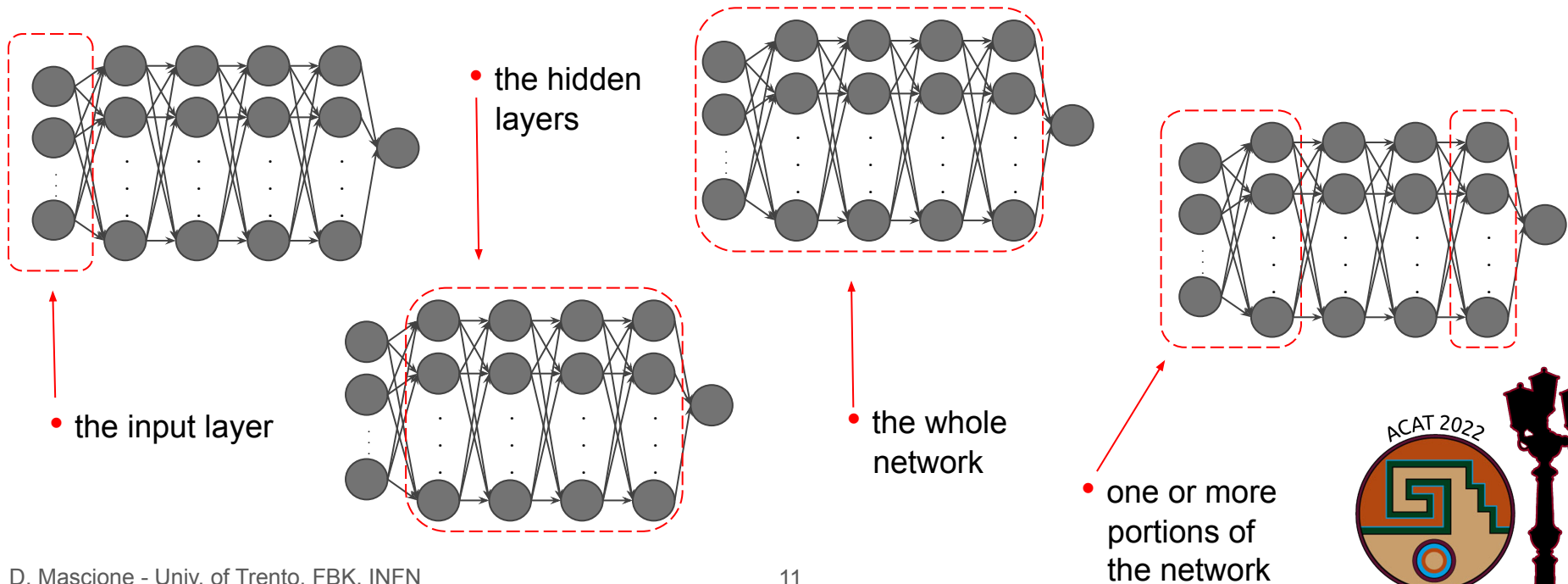- it can determine the most suitable **network architecture**

AutoPruner

# AutoPruner: a different pruning strategy

- it can prune **nodes**

- it prunes **during training**

- the number of nodes to be pruned can be determined by the **user**

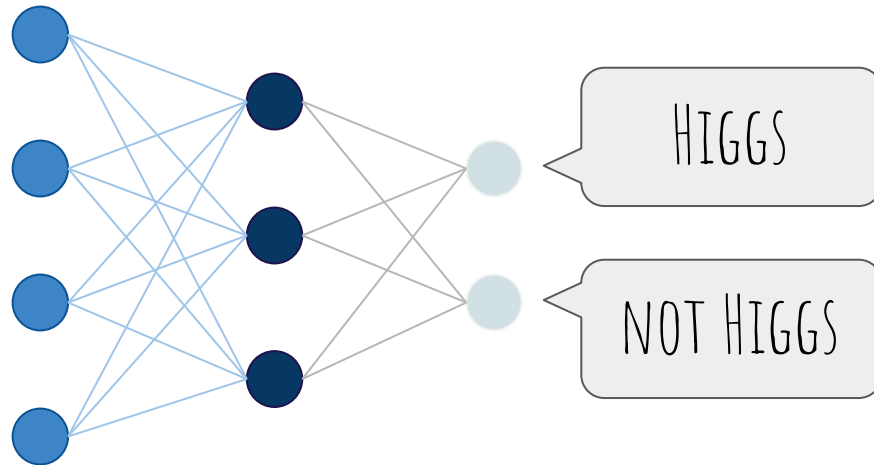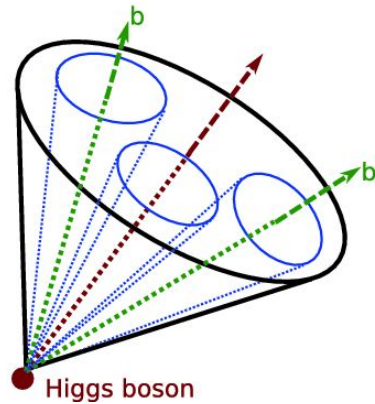- it can determine the most suitable **network architecture**

AutoPruner

# AutoPruner: a different pruning strategy

- it can prune **nodes**

- it prunes **during training**

- the number of nodes to be pruned can be determined by the **user**

- it can determine the most suitable **network architecture**

# AutoPruner: a different pruning strategy

- it can prune **nodes**

- it prunes **during training**

- the number of nodes to be pruned can be determined by the **user**

- it can determine the most suitable **network architecture**

# Pruning with AutoPruner

With AutoPruner you can **choose** which part of the network you want to prune

- the hidden layers

- the input layer

- the whole network

- one or more portions of the network

ACAT 2022
BARI

# Use case

Identify jets that contain both the *b* quarks from boosted Higgs decay in *pp* collision experiments using Deep Neural Networks
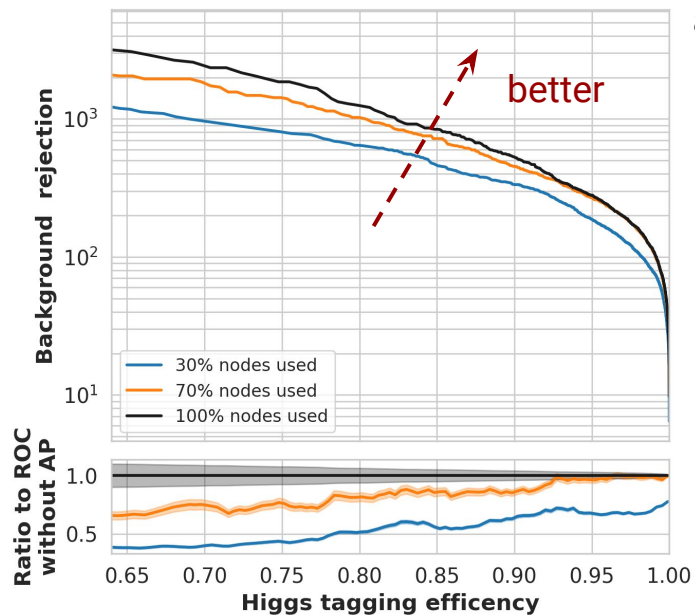
# *Results*

The performance increases with the percentage of nodes used, as expected: AutoPruner is really **switching off** nodes
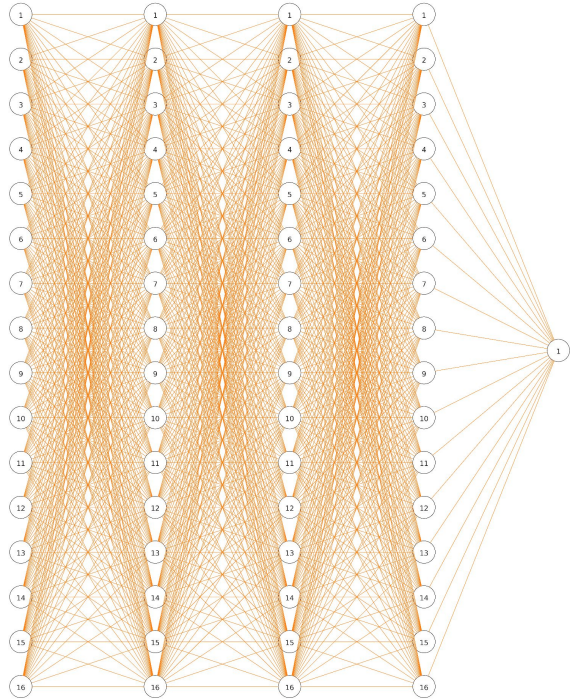


better

ACAT 2022
BARI

# Results

The performance increases with the percentage of nodes used, as expected: AutoPruner is really **switching off** nodes
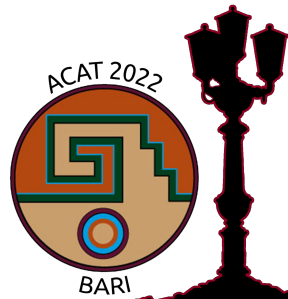


better

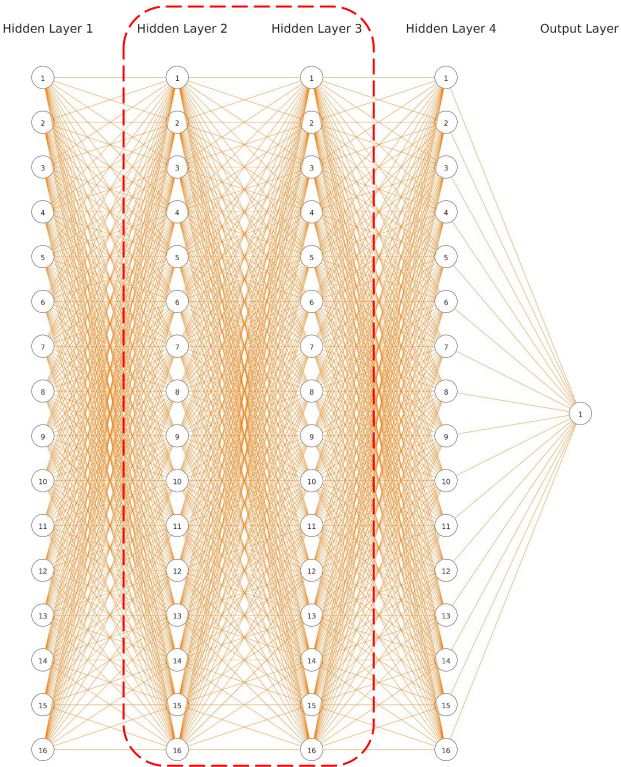The total number of nodes used is **always** equal to the required number

ACAT 2022
BARI

# Models' comparison



Hidden Layer 1    Hidden Layer 2    Hidden Layer 3    Hidden Layer 4    Output Layer

ACAT 2022
BARI

# *Models' comparison*



Hidden Layer 1   Hidden Layer 2   Hidden Layer 3   Hidden Layer 4   Output Layer

ACAT 2022

BARI

# Models' comparison

# Models' comparison

# Models' comparison

# AutoPruner in Convolutional Neural Networks



AutoPruner can be used with CNNs to prune filters during training
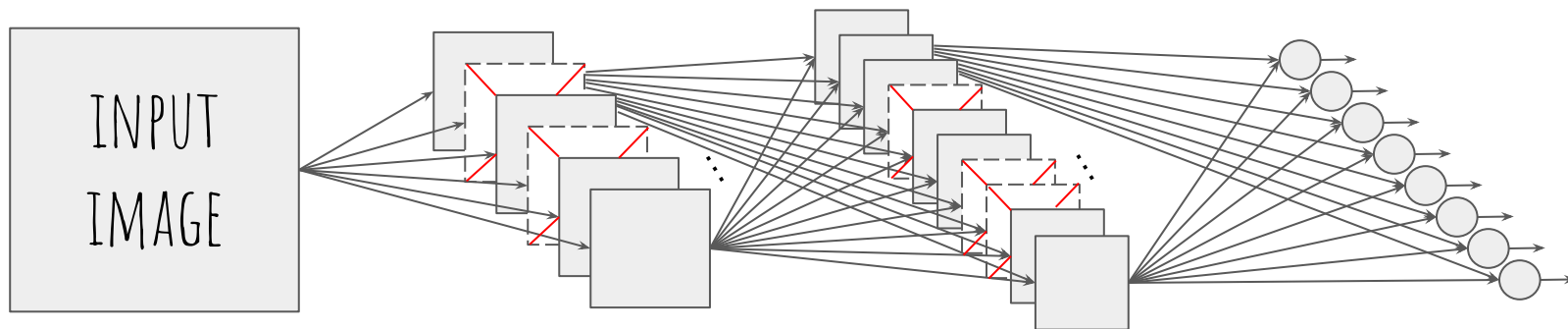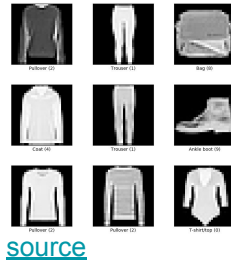
# AutoPruner in Convolutional Neural Networks



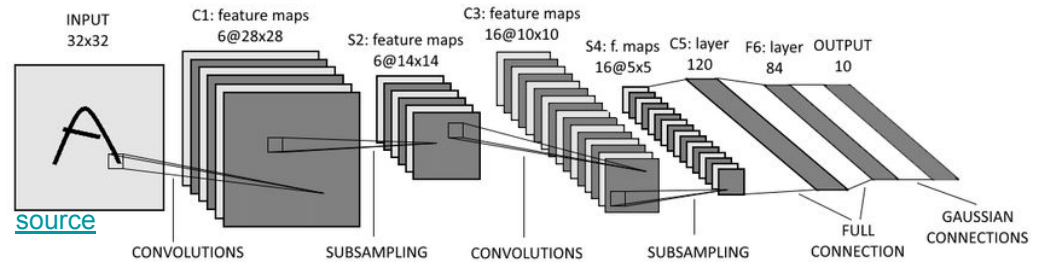AutoPruner can be used with CNNs to prune filters during training

# Preliminary results

Use case:
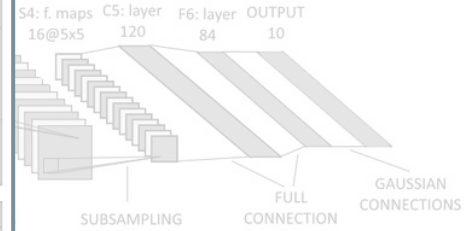


source

+



source

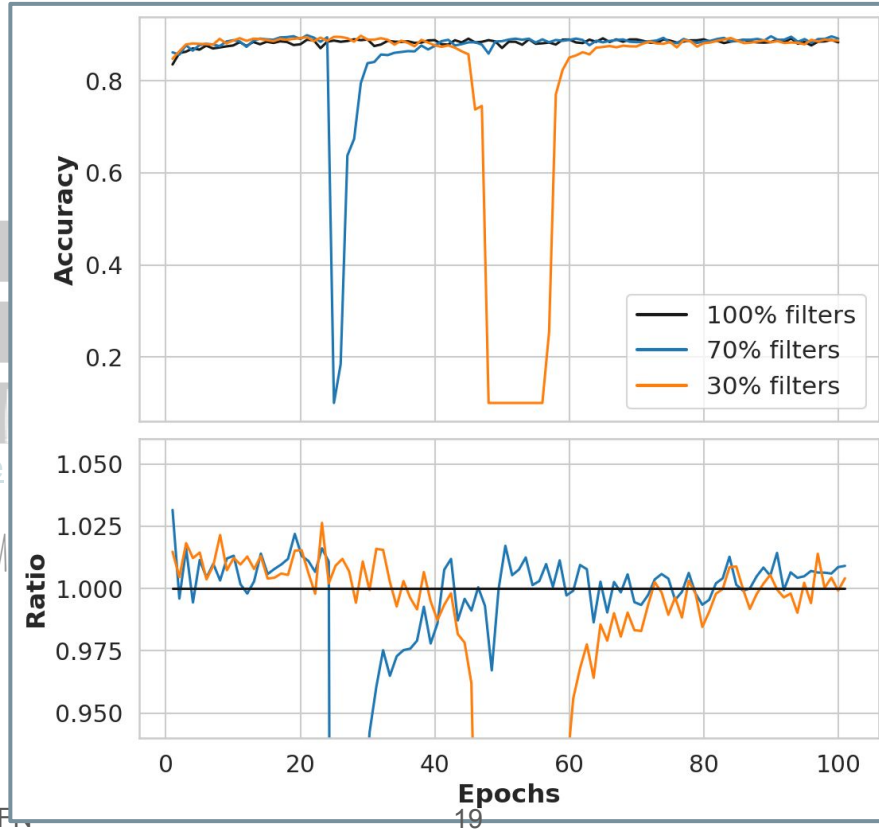Fashion-MNIST dataset

LeNet-5 architecture

# Preliminary results



Use case:

source

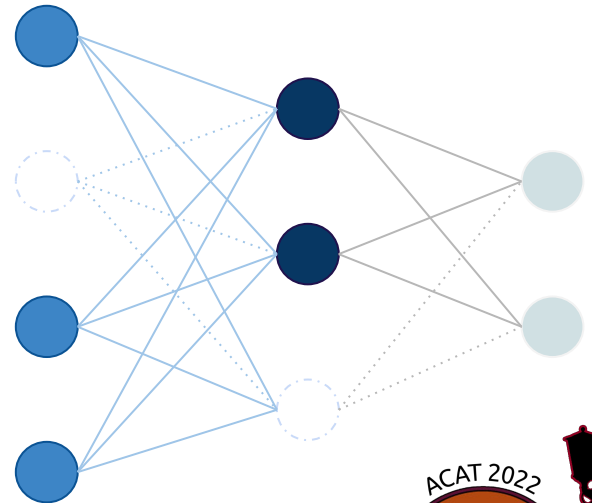Fashion-M...                                                ...cture

# *Conclusions*

We introduced the AutoPruner approach to **effectively prune** Deep Neural Networks during training.

AutoPruner proved to be:

- **simple** to incorporate
- **effective** and **successful** in reducing the networks' size
- **fast** (pruning during training, no need to fine tune)
- very **understandable**

Further developments are focusing on:

- apply AutoPruner to Convolutional Neural Networks
- investigate feature selection with AutoPruner

# Thanks!

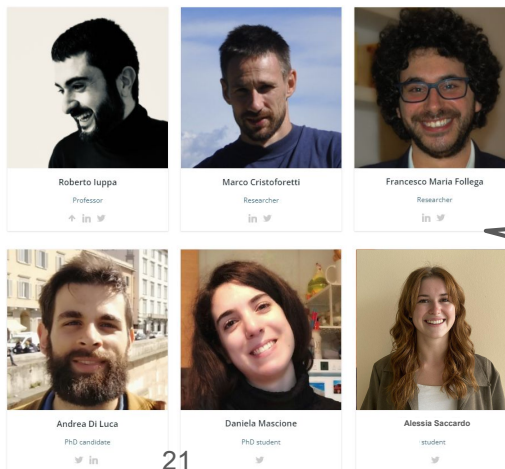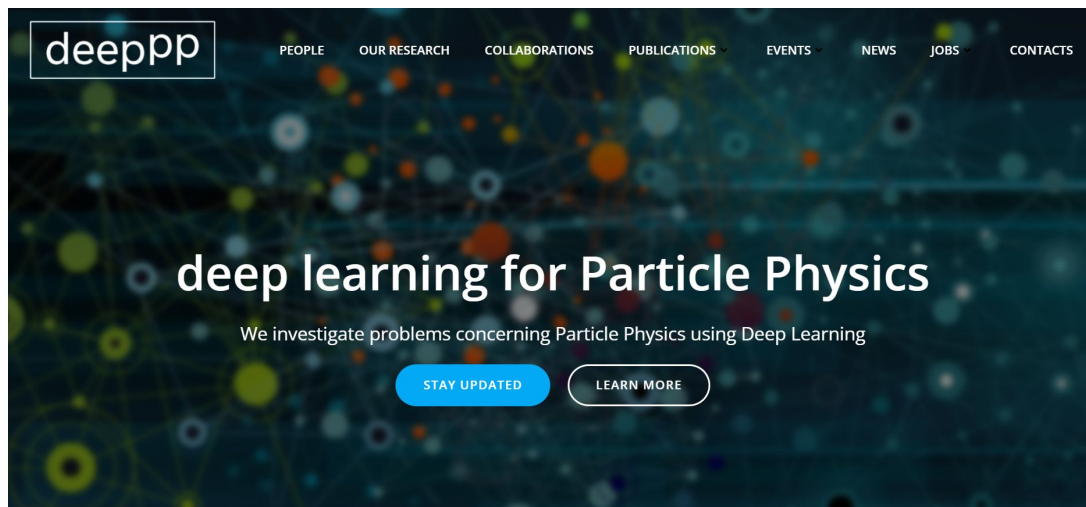Want to know more about Deep Learning applications in Particle Physics?

Awesome!

Visit

https://www.deeppp.eu/

deeppp

D. Mascione - Univ. of Trento, FBK, INFN



21