# Convolutional layer & GPUs
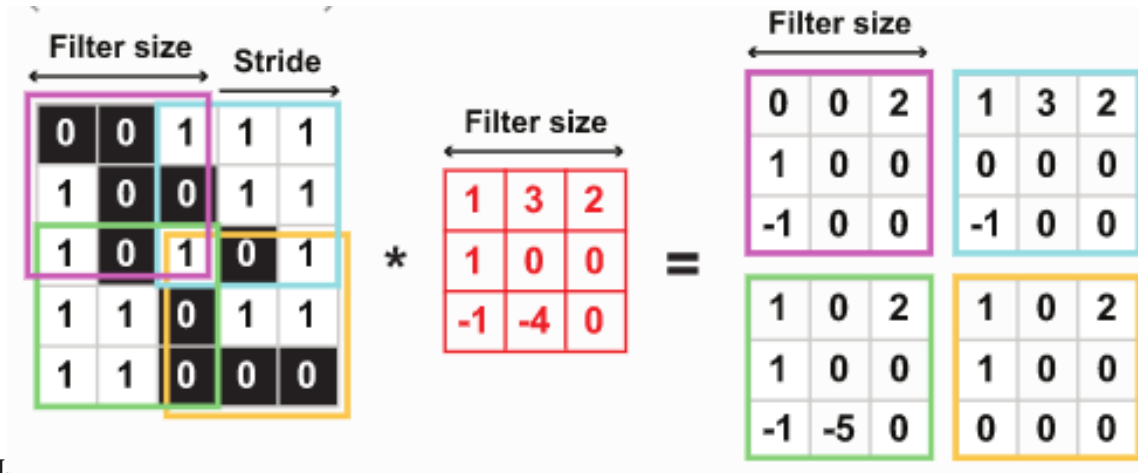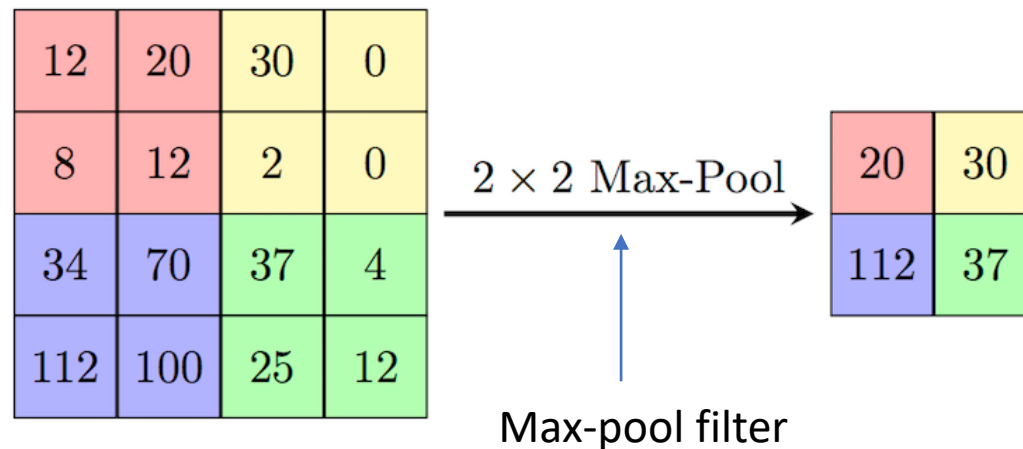
# Convolutional layer

- The convolutional layer prepares the input data for the fully connected neurons (hidden layers)

- Procedures happening the convolutional layer:
  - Filtering
  - Pooling
  - Flatten & D

- Filters (kernel):
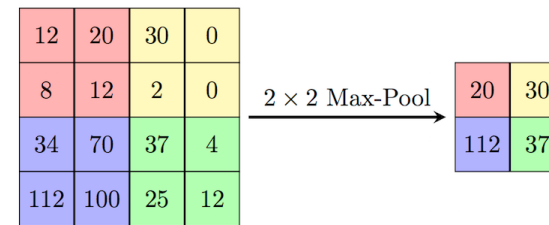  - In charge to detect features
    - In images:
      - Can detect for example, lines, corners
      - Can also be specific: eyes, hair, etc.
  - Is a matrix of m * n. Size of the matrix depends on the application
    - Initialized randomly
    - Stride: how many "cells" the filter moves for a dot product
  - Operation:
    - A dot product operation is executed between the input data and the filter

- Pooling:
  - To reduce dimensionality
  - Eliminates noise from data

  - Max pooling: after the dot product of the filter and input is executed
    - A max pooling filter (of m*n) selects the max value and creates a new max pool.



Max-pool filter

- Flatten & Dense:
  - Flatten: For the neurons to perform their own operations, the input data needs to be an appropriate format:



The output will be the classification
based on the features extracted
through the previous layers
and their different filters
(classification: outputs are
mapped to labels)

The vector (tensor) is passed as input of the FC

The output if the max-pool is a matrix
Can be of n-dimensions

**Flattening**: "Converts" n-dimension into a "readable" <u>vector (tensor)</u>

- Dense:
  - Method to create a fully connected neural network
  - Specifies the number of neurons per layer and activation method



Input     Convolutional     FC     FC     Output

Pooling/Filtering ⟶ Flatten/Dense

- How is a CNN processed when GPUs are available?

# GPU ARCHITECTURE

Streaming Multiprocessor (SM)

Many CUDA Cores per SM

Architecture dependent

H100 SM has *128 cores*

Special-function units

cos/sin/tan, etc.

Shared mem + L1 cache

Thousands of 32-bit registers

H100 PCIe has a total of *14,592 cores*

- How is a CNN processed when GPUs are available?

**PROCESSING FLOW**



1. Copy input data from CPU memory to GPU memory
2. Load GPU program and execute, caching data on chip for performance

- How is a CNN processed when GPUs are available?

## PROCESSING FLOW



1. Copy input data from CPU memory to GPU memory
2. Load GPU program and execute, caching data on chip for performance
3. Copy results from GPU memory to CPU memory

- How is a CNN processed when GPUs are available?



**5 WAYS TO ACCELERATE WITH GPUS**

| Applications | Libraries | OpenACC Directives | CUDA Programming | Standard Language Parallelism |
|---|---|---|---|---|
| Get straight to the science! | "Drop-in" Acceleration | Easily Accelerate Applications | Maximum Performance | Maximum Flexibility |

Flexibility →

← Accessibility

THE UNIVERSITY OF SOUTHERN MISSISSIPPI.

- How is a CNN processed when GPUs are available?

**ARTIFICIAL INTELLIGENCE**
- PyTorch
- MXNet
- TensorFlow

...

**CLIMATE & WEATHER**
- Cosmos
- Gales
- WRF

...

**COMPUTATIONAL FINANCE**
- O-Quant Options Pricing
- MUREX
- MISYS

...

**DATA SCIENCE & ANALYTICS**
- Anaconda
- H20
- OmniSci

...

**FEDERAL DEFENSE & OTHER**
- ArcGIS Pro
- EVNI
- SocetGXP
- Cyllance
- FaceControl

...

**LIFE SCIENCES**
- Amber
- LAMMPS
- GROMACS
- NAMD
- Relion
- VASP

...

**MANUFACTURING, CAD, & CAE**
- Ansys Fluent
- Abaqus SIMULIA
- AutoCAD
- CST Studio Suite

...

**MEDIA & ENTERTAINMENT**
- DaVinci Resolve
- Premiere Pro CC
- Redshift Renderer

...

**MEDICAL IMAGING**
- aidoc
- PowerGrid
- RadiAnt

...

**OIL & GAS**
- Echelon
- RTM
- SPECFEM3D

...

**RETAIL**
- Everseen
- Deep North
- Third Eye Labs
- AWM
- Malong
- Clarifai
- Antuit

...

**SUPERCOMPUTING & HER**
- Chroma
- GTC
- MILC
- QUDA
- XGC

...

- How is a CNN processed when GPUs are available?



**ALGORITHMS**
GPU-accelerated Scikit-Learn

| Category | Algorithms |
|---|---|
| Classification / Regression | Decision Trees / Random Forests<br>Linear/Lasso/Ridge/ElasticNet Regression<br>Logistic Regression<br>K-Nearest Neighbors<br>Support Vector Machine Classification and Regression<br>Naive Bayes |
| Inference | Random Forest / GBDT Inference (FIL) |
| Preprocessing | Text vectorization (TF-IDF / Count)<br>Target Encoding<br>Cross-validation / splitting |
| Clustering Decomposition & Dimensionality Reduction | K-Means<br>DBSCAN<br>Spectral Clustering<br>Principal Components<br>Singular Value Decomposition<br>UMAP<br>Spectral Embedding<br>T-SNE |
| Time Series | Holt-Winters<br>Seasonal ARIMA / Auto ARIMA |

Cross Validation

Hyper-parameter Tuning

More to come!