

Run: 204153
Event: 35369265
2012-05-30 20:31:28 CEST



ALGORITMI VELOCI DI DEEP LEARNING SU PROCESSORI FPGA PER IL TRIGGER DI MUONI DI LO NELLA FASE II DELL'ESPERIMENTO ATLAS A LHC

IFAE2019 - NAPOLI
9 APRILE 2019

LUIGI SABETTA

MOTIVAZIONI

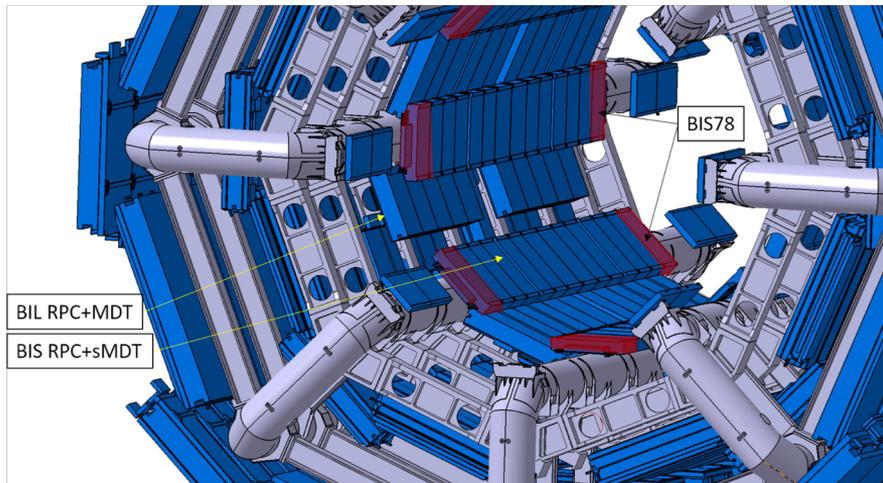
IL TRIGGER MUONICO DI LIVELLO 0 DI ATLAS SUBIRÁ UN UPGRADE COMPLETO PER HL-LHC

PARAMETRI DELLA MACCHINA:

- MAGGIORE PILE UP: 30-40 → 200
- LUMINOSITÀ: (2 → 7.5) × 10³⁴ cm⁻²s⁻¹

MAGGIOR RATE DI HIT NEL SETTORE MUONICO:

FINO A 600 HZ/CM²



POTENZIAMENTO DI ATLAS:

NUOVO PROCESSORE DI TRIGGER

- SISTEMA BASATO SU FPGA
 - VIRTEX ULTRASCALE+ XCVU13P
 - CELLE LOGICHE (K): 3780
 - MEMORIA (MB): 455
 - GTY TRANSRECEIVERS (32.75 GB/s): 128
 - I/O PINS: 832

NUOVA STAZIONE DI TRIGGER

- NUOVI LAYER RPC

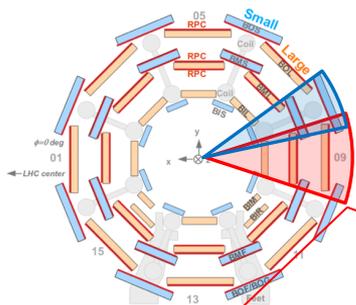
ALGORITMO DI TRIGGER MIGLIORATO

- NECESSITÀ DI ESSERE MOLTO FLESSIBILE E VELOCE

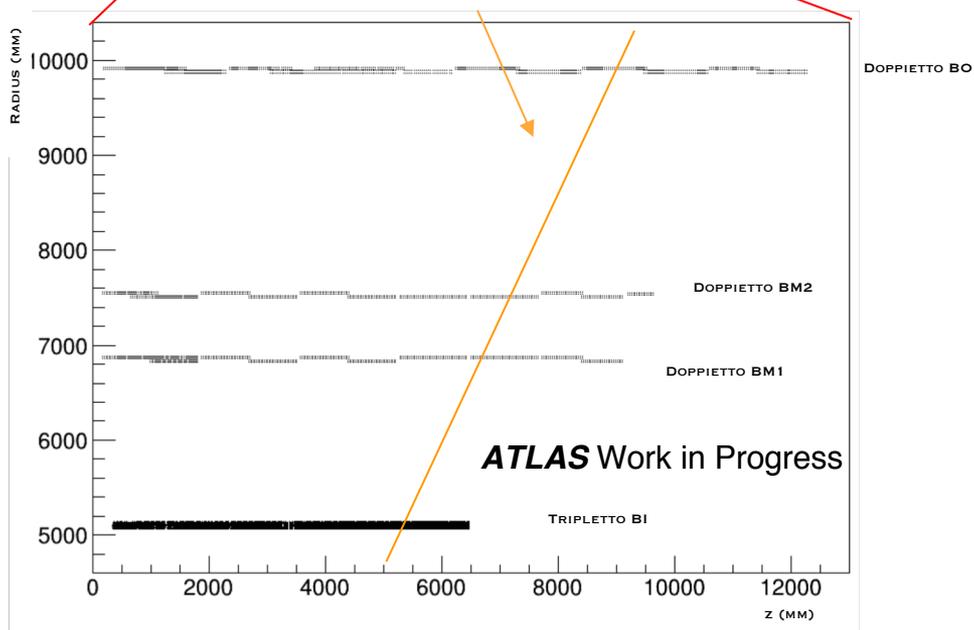
RETI NEURALI → VALIDO CANDIDATO

[ATLAS-TDR-026](#)

INFORMAZIONI DISPONIBILI PER IL TRIGGER



MUONE DI IMPULSO TRASVERSO INFINITO

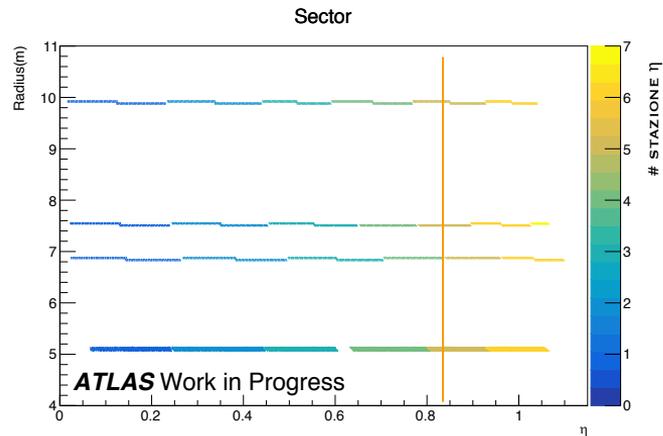


RESISTIVE PLATE CHAMBER(RPC):

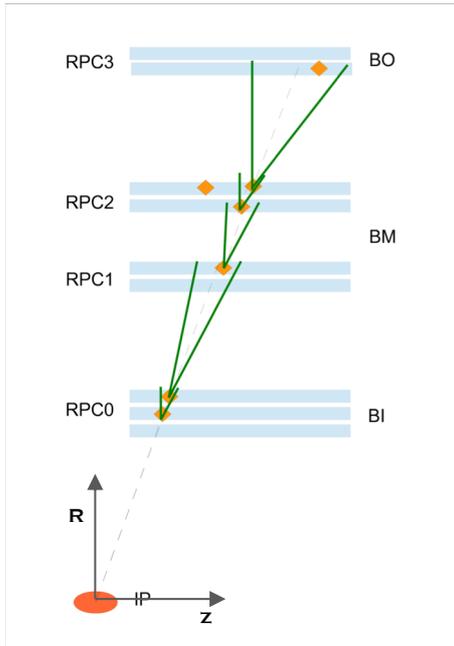
- DETECTOR A GAS VELOCI
- SORGENTE DI TRIGGER IN ATLAS

SETTORE:

- SI DIVIDE ATLAS IN 2 PARTI ($z \geq 0$)
- SI DIVIDE LA CORDINATA ϕ IN 8 SETTORI, CIASCUNO DEI QUALI VIENE A SUA VOLTA DIVISO IN 2 (SMALL E LARGE)
- SI PRENDONO IN CONSIDERAZIONE TUTTE LE STRIP DEGLI RPC DEL SETTORE OTTENUTO



STRATEGIA STANDARD DI TRIGGER



[ATLAS-TDR-026](#)

ALGORITMO STANDARD:

- CONTROLLA LA PRESENZA DI COINCIDENZE ALL'INTERNO DI FINESTRE APERTE SEQUENZIALMENTE IN BASE ALLA PRESENZA DI SEGNALE NEI LAYERS PRECEDENTI

👍 STABILE ED AFFIDABILE

👍 BUONE PERFORMANCE

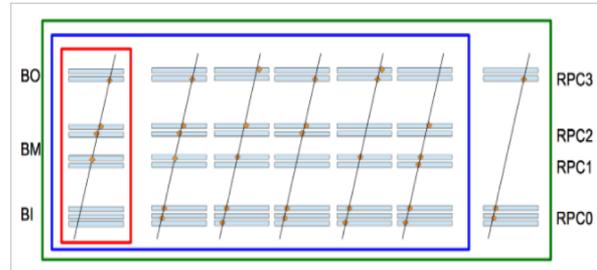
👎 LE FINESTRE DI COINCIDENZA DEVONO ESSERE REGOLATE “A MANO”

- LA DIMENSIONE DELLE FINESTRE DIPENDE DALLA SOGLIA DI TRIGGER IN IMPULSO TRASVERSO (p_T)
- FORTE DIPENDENZA DALLE MODALITÀ DI TRIGGER E DALLA GEOMETRIA LOCALE DEL RIVELATORE

👎 ASSUME TRACCE CHE PUNTANO AL VERTICE PRIMARIO DI INTERAZIONE

- VERTICI DI DECADIMENTO SECONDARI LONTANI?

👎 NON FORNISCE NESSUN TIPO DI STIMA SUL p_T



TRIGGER E NEURAL NETS (NN)

OUTPUT MINIMALE TRIGGER → FINO A TRE CANDIDATI

UNA NN PUÓ ESSERE SFRUTTATA PER:

- INDIVIDUARE UN CANDIDATO MUONE (=ALMENO UN MUONE CON $p_T > \text{SOGLIA}$)
- EFFETTUARE REGRESSIONE SU p_T ED η DEL MUONE CON p_T PIÚ ALTO (LEADING)
- EFFETTUARE REGRESSIONE SU p_T ED η DEL SECONDO MUONE CON p_T PIÚ ALTO (SUBLEADING) → **BONUS!**
- CONTROLLARE (CLASSIFICARE) IL NUMERO DI MUONI → **BONUS!**

REGRESSIONE SUL p_T A LIVELLO 0:

- L'ALGORITMO STANDARD RISPONDE SEMPLICEMENTE ALLA DOMANDA “È SOPRA SOGLIA?”
(IMPOSTANDO DIVERSE SOGLIE SI PUÒ INDIVIDUARE UN RANGE DI IMPULSO)

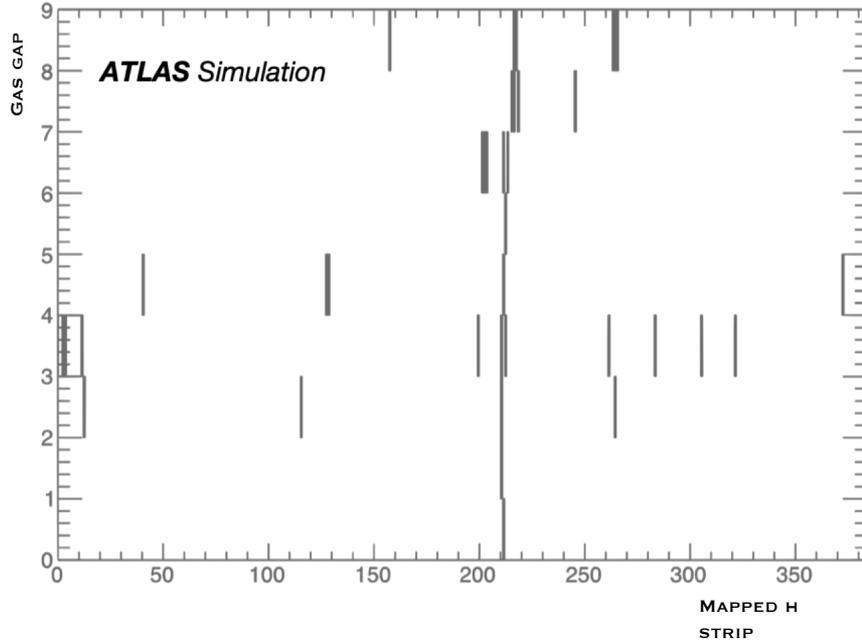
LA NN VIENE ADDESTRATA PER FORNIRE UN OUTPUT 5D

$(p_T^{\text{lead}} \quad \eta^{\text{lead}} \quad p_T^{\text{sublead}} \quad \eta^{\text{sublead}} \quad N^{\text{muons}})$

0 → NESSUN MUONE
1/2 → UNO/DUE MUONI
3 → PIÚ DI 2 MUONI

DA MAPPA DI STRIP A IMMAGINI

[HTTPS://TWIKI.CERN.CH/TWIKI/BIN/VIEW/ATLASPUBLIC/L1MUONTRIGGERPUBLICRESULTS](https://twiki.cern.ch/twiki/bin/view/ATLASPUBLIC/L1MUONTRIGGERPUBLICRESULTS)



ESEMPIO DI UN MUONE CON $P_T=19$ GEV + NOISE

MAPPAGGIO $\rightarrow i^{\text{strip}} = \frac{\eta_{\text{hit}} - \eta_{\text{min}}}{\eta_{\text{max}} - \eta_{\text{min}}} \times 384$

SI POSSONO COSÍ COSTRUIRE IMMAGINI DA USARE
COME INPUT PER UNA NN:

- NEL CASO MINIMALE 384 X 9 PIXEL
- MUONE \rightarrow RETTA

VOGLIAMO FARE REGRESSIONE SUL p_T DEL MUONE

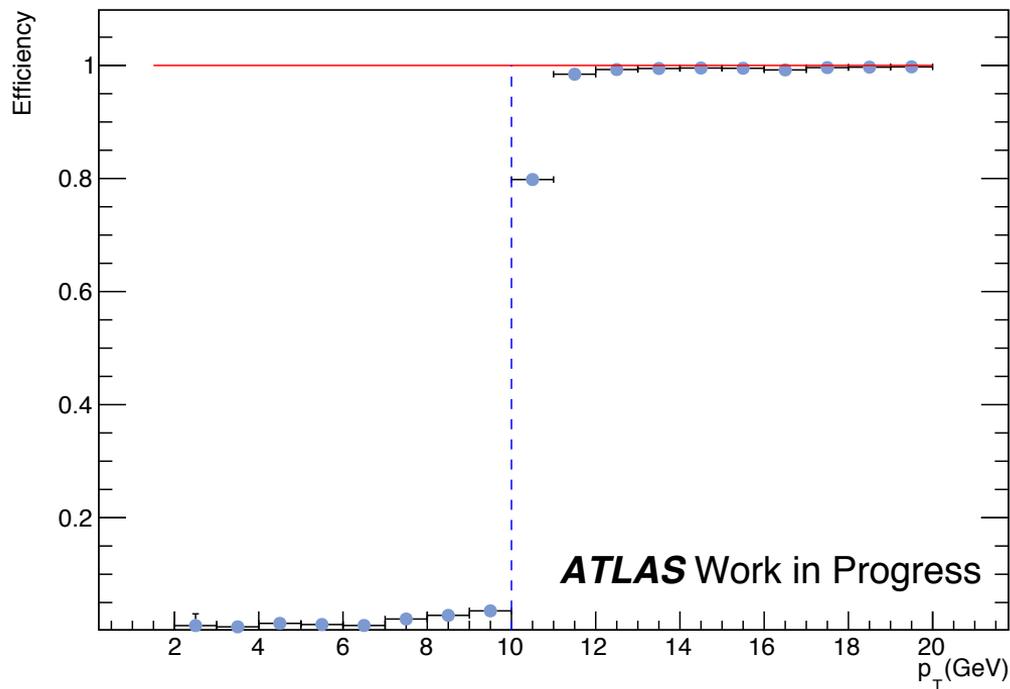
PROBLEMA:

- LE IMMAGINI SONO ALTAMENTE
“SPARSIFICATE”

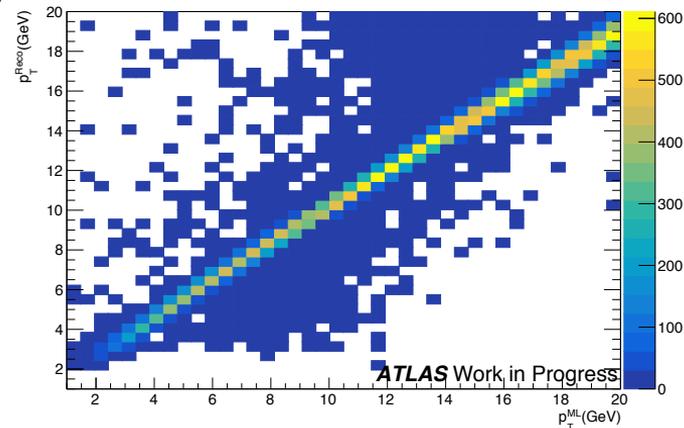
RETE NEURALE DENSA - PERFORMANCE

RETE NEURALE COMPLETAMENTE CONNESSA (2 DA LAYERS DA 2048 NEURONI)

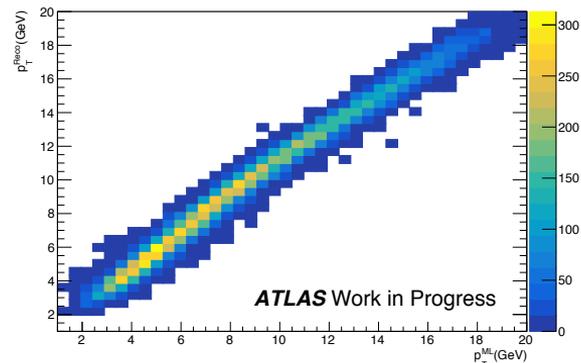
INPUT → IMMAGINI LINEARIZZATE



~11M DI PARAMETRI → TROPPO GRANDE DA SINTETIZZARE

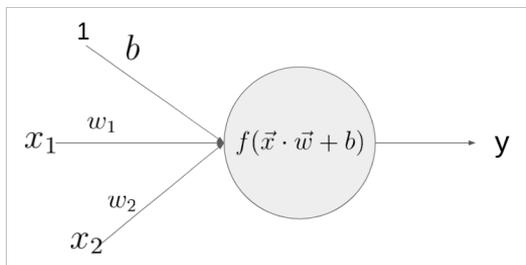


MUONE LEADING (↑) E SUBLEADING (↓)



RETI CONVOLUZIONALI

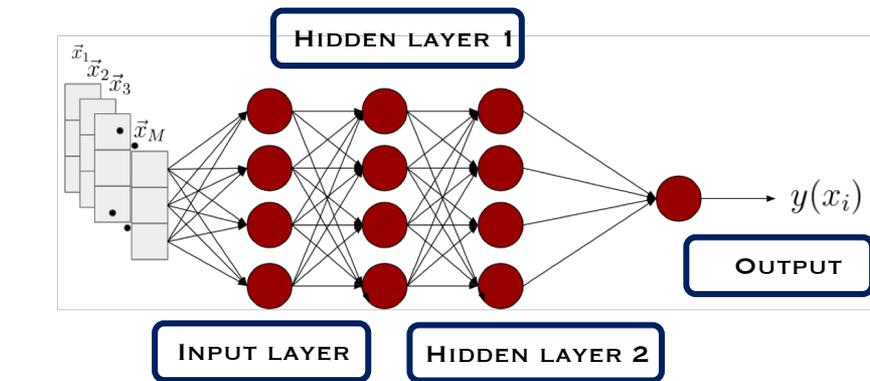
NEURONE



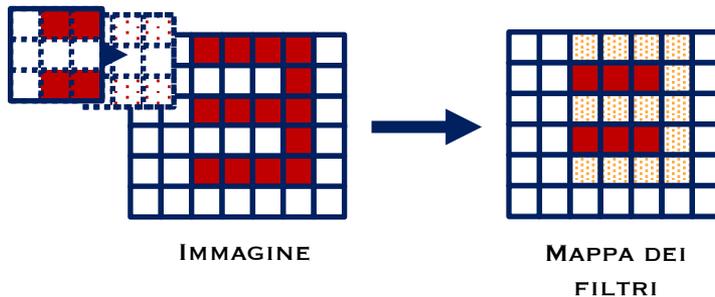
RETE NEURALE CONVOLUZIONALE (CNN):
RETI NEURALI PROFONDE
OTTIMIZZATE PER L'IDENTIFICAZIONE
DI IMMAGINI.



EFFICACI PER PROBLEMI CHE
PRESENTANO SIMMETRIE
TRASLAZIONALI E ROTAZIONALI

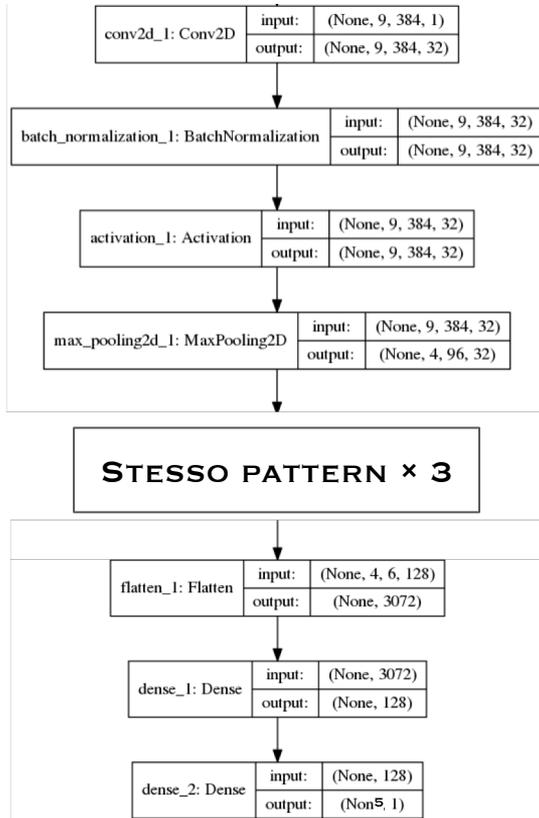


FILTRO



IN QUESTO MODO IL NUMERO DI PARAMETRI SI RIDUCE NOTEVOLMENTE

STRUTTURA DELLA CNN FLOATING POINT



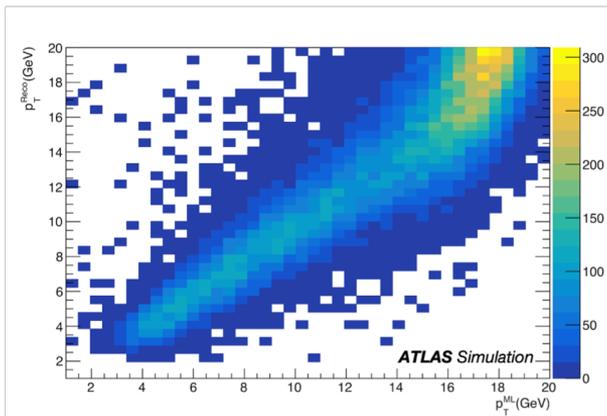
ARCHITETTURA:

- (CONV2D + BATCH NORM. + MAX POOLING) X 3
- (DENSE) X 2 LAYERS PER ARRIVARE ALL'OUTPUT (5D)

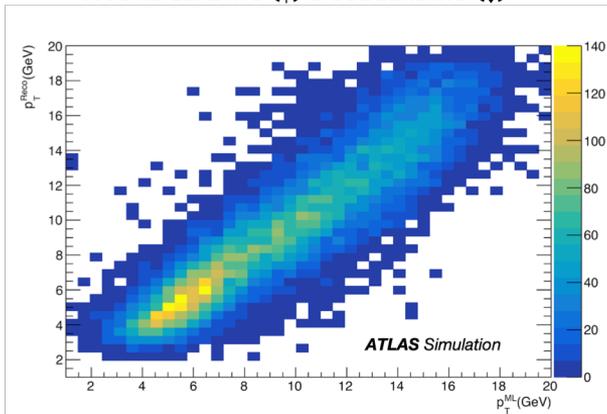
NUMERO TOTALE DI PARAMETRI: 500K

- MENO DI UNA NN DENSE PURA

PERFORMANCE DELLA CNN FP- QUANTITÀ FISICHE

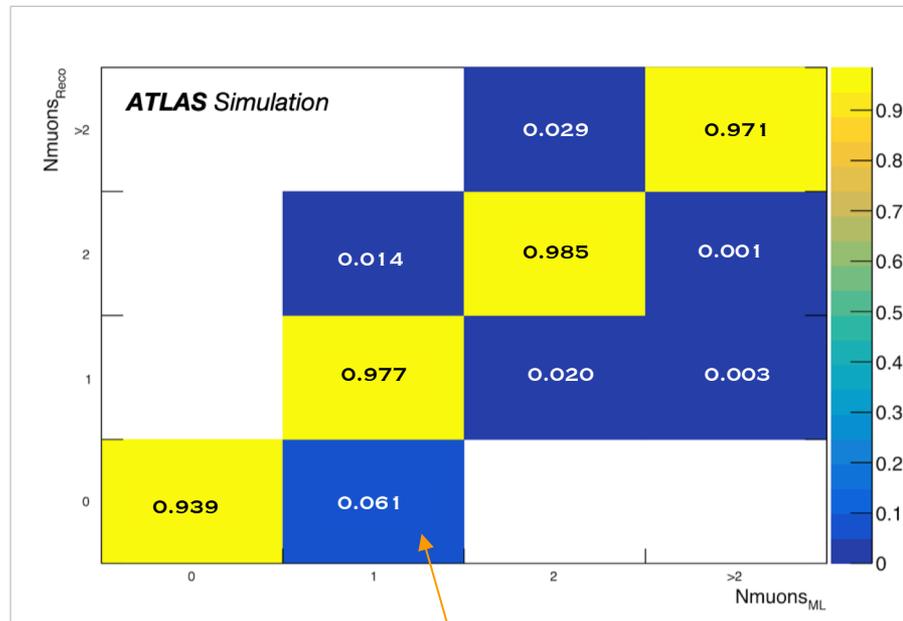


MUONE LEADING (\uparrow) E SUBLEADING (\downarrow)



LE QUANTITÀ FISICHE INTERESSANTI SONO BEN RAPPRESENTATE

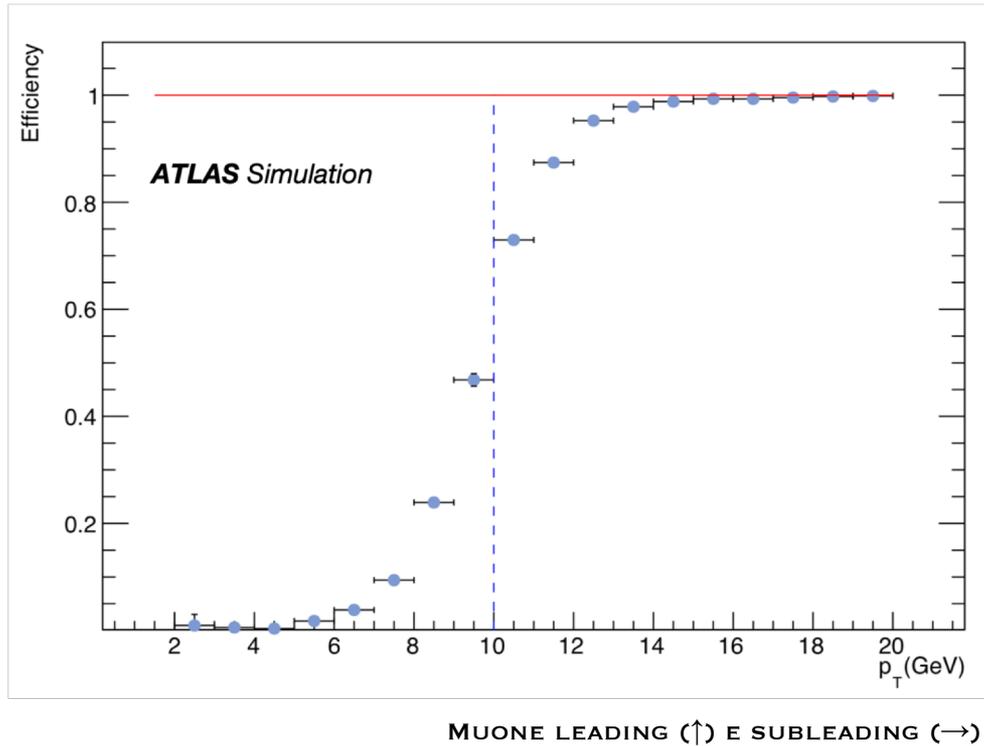
GLI EVENTI RISULTANO BEN CLASSIFICATI PER QUANTO RIGUARDA n^{muoni}



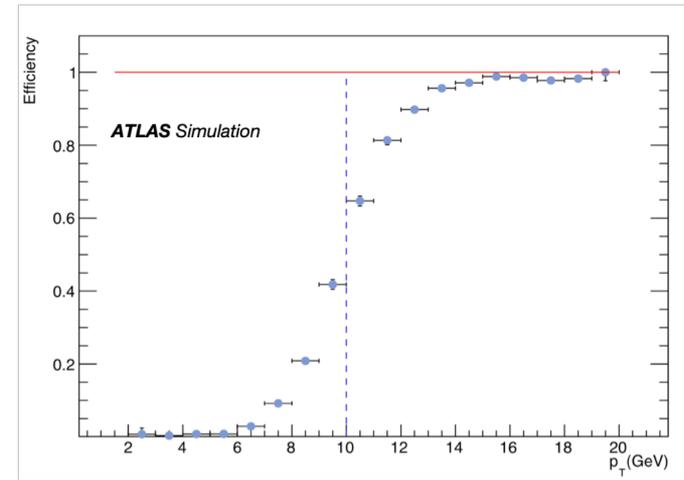
QUASI TOTALMENTE MUONI CON p_T FRA 0-2 GeV: RIMOSI DAL TRIGGER

[HTTPS://TWIKI.CERN.CH/TWIKI/BIN/VIEW/ATLASPUBLIC/L1MUONTRIGGERPUBLICRESULTS](https://twiki.cern.ch/twiki/bin/view/ATLASPUBLIC/L1MUONTRIGGERPUBLICRESULTS)

PERFORMANCE DELLA CNN FP- TRIGGER



LE CURVE DI EFFICIENZA DI TRIGGER SONO SODDISFACENTI SIA PER IL MUONE LEADING CHE PER IL SUB-LEADING



IMPLEMENTAZIONE - CNN TERNARIA

LA SECTOR LOGIC DEL TRIGGER DI LIVELLO 0 OPERERÁ SU FPGA

É NECESSARIO QUINDI SINTETIZZARE QUESTA CNN SU DI UNA FPGA

- NN STANDARD FUNZIONANO CON PESI DESCRITTI DA NUMERI CON PRECISION FLOATING POINT A 32 BIT
- PESI FP NON SONO LA SCELTA OTTIMALE
 - GROSSO IMPIEGO DI RISORSE LOGICHE

SI PUÓ PENSARE DI REALIZZARE UNA CNN TERNARIA:

- PESI =-1, 0, +1 (DESCRITTI DA SOLI DUE BIT DI MEMORIA)

UNA NN TERNARIA DETERMINA UNA PERDITA DI PRECISIONE

- POCHI PUNTI PERCENTUALI RISPETTO AD UNA NN FP32 CON LA STESSA STRUTTURA

MA DIMENSIONI MINORI

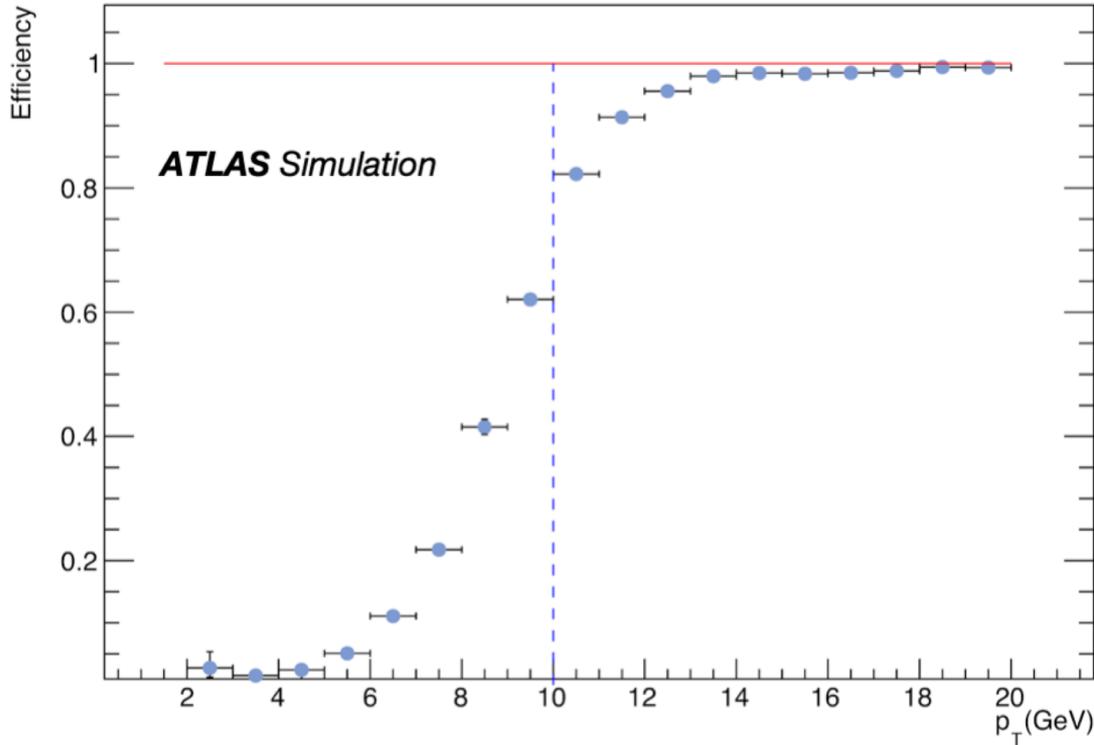
- INFERIORE CONSUMO DI RISORSE LOGICHE
- FINO A 16 VOLTE PIÚ PICCOLA RISPETTO AD UNA NN FP32



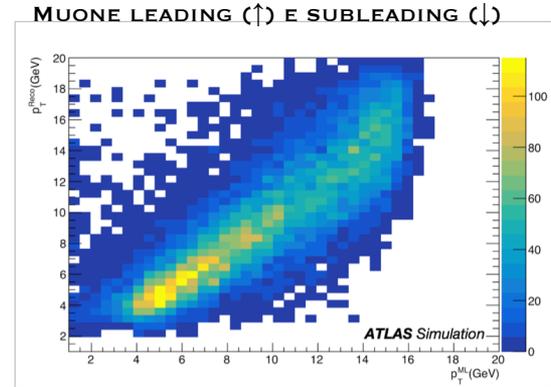
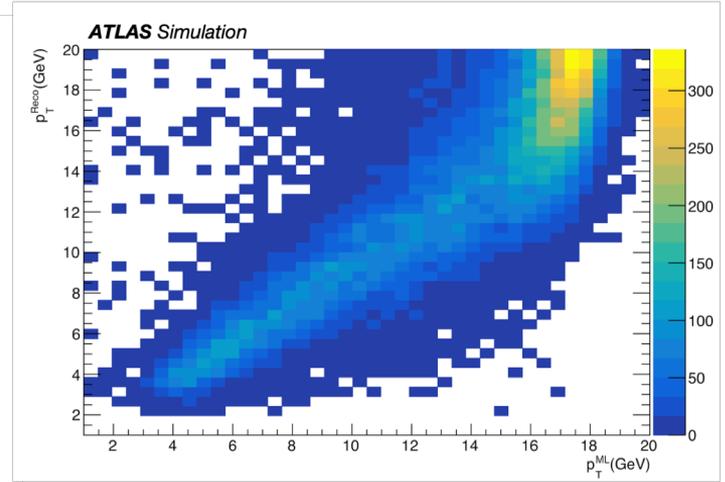
PUÓ IN LINEA DI PRINCIPIO ESSERE RESA PIÚ PROFONDA

- PIÚ LAYER RECUPERANO LA PRECISIONE PERSA

CNN TERNARIA- PERFORMANCE



<https://twiki.cern.ch/twiki/bin/view/ATLASPUBLIC/L1MUONTRIGGERPUBLICRESULTS>



STATO ATTUALE DEI LAVORI

UNA CNN PUÓ ESSERE UTILIZZATA PER IMPLEMENTARE L'ALGORITMO DI TRIGGER MUONICO DI LIVELLO 0 PER L'UPGRADE DI FASE II DI ATLAS

- BUONA RISOLUZIONE E RISPOSTA
- REGRESSIONE SU p_T ED η DI MUONE LEAD. E SUBLEAD.

ABBIAMO UNA RETE TERNARIA CONVOLUZIONALE CHE MIGLIORA LE PRESTAZIONI DELL'ALGORITMO CONVENZIONALE

- IN TERMINI DI RISORSE LOGICHE, PUÓ ESSERE SINTETIZZATA NELLA FPGA CHE SI PIANIFICA DI IMPIEGARE PER IL TRIGGER LO

ABBIAMO SINTETIZZATO IN UNA FPGA CON SUCCESSO PER IL MOMENTO SEMPLICI ARCHITETTURE DENSE E CONV2D

- MODELLI FIRMWARE REALIZZATI USANDO HLS4ML

ATTUALMENTE STIAMO LAVORANDO PER RIUSCIRE A SINTETIZZARE IL NOSTRO MODELLO DI RETE TERNARIA CONVOLUZIONALE

NN DENSE F32 CON 3 LAYER

* Summary:					
Name	BRAM_18K	DSP48E	FF	LUT	URAM
DSP	-	-	-	-	-
Expression	-	-	0	6	-
FIFO	-	-	-	-	-
Instance	13	3376	44226	131419	-
Memory	-	-	-	-	-
Multiplexer	-	-	-	36	-
Register	-	-	4450	-	-
Total	13	3376	48676	131461	0
Available	960	1824	433920	216960	64
Utilization (%)	1	185	11	60	0

NN DENSE TERNARIA CON 3 LAYER

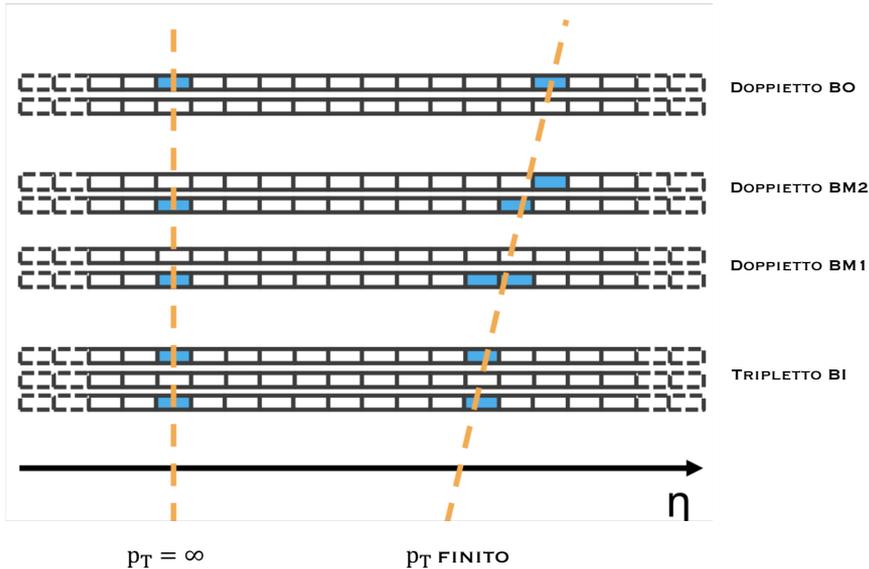
* Summary:					
Name	BRAM_18K	DSP48E	FF	LUT	URAM
DSP	-	-	-	-	-
Expression	-	-	0	6	-
FIFO	-	-	-	-	-
Instance	-	123	9626	59344	-
Memory	-	-	-	-	-
Multiplexer	-	-	-	36	-
Register	-	-	5906	-	-
Total	0	123	15532	59386	0
Available	960	1824	433920	216960	64
Utilization (%)	0	6	3	27	0

GRAZIE PER L'ATTENZIONE

BACKUP

MAPPAGGIO IN STRIP- η VS LAYERS

$$i^{\text{strip}} = \frac{\eta_{\text{hit}} - \eta^{\text{min}}}{\eta^{\text{max}} - \eta^{\text{min}}} \times 384$$



IL MAPPAGGIO PREVIENE ALCUNE DIPENDENZE GEOMETRICHE, COME:

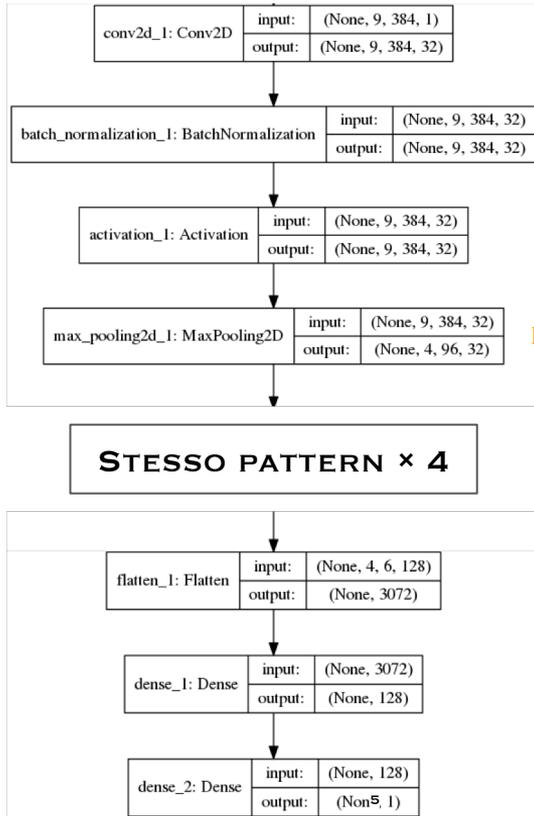
- **BUCHI**
- **SOVRAPPOSIZIONE DI CAMERE**

DA NOTARE CHE UN MUONE DI IMPULSO TRASVERSO INFINITO RESTA SEMPRE UNA LINEA VERTICALE, INDIPENDENTEMENTE DA η

IL MAPPAGGIO IN STRIP FITTIZIE PERMETTE DI AUMENTARE/DIMINUIRE LA GRANULARITÀ DEL DETECTOR A PIACERE:

- **NEL CASO MINIMALE, AD OGNI STRIP REALE CORRISPONDE UNA STRIP FITTIZIA**

STRUTTURA DELLA CNN TERNARIA



FILTRI [3 × 3]

POOL. [2 × 4]

ARCHITETTURA:

- (CONV2D + BATCH NORM. + MAX POOLING) x4
- STESSA STRUTTURA DELLA RETE FP32
- (DENSE)x2 LAYERS PER ARRIVARE ALL'OUTPUT (5D)

NUMERO TOTALE DI PARAMETRI: 1M

- MAGGIORE RISPETTO ALLA CNN FP
- MA QUI OGNI PESO PUÓ ESSERE SOLAMENTE +1/0/-1
- POTENZIALMENTE SOLO 2 BIT

PER MASSIMIZZARE LE PRESTAZIONI DELLA RETE, IL TRAINING (OPERATO PRELIMINARMENTE SU GPU) VIENE EFFETTUATO IN MANIERA NON TERNARIA

NN SU FPGA – COME IMPLEMENTARLE



ARCHITETTURA
IN PYTHON

PER IMPLEMENTARE UNA NN IN UNA FPGA BISOGNA
TRADURRE ARCHITETTURE TRADIZIONALI REALIZZATE CON
PACCHETTI PER LA COSTRUZIONE DI RETI NEURALI (KERAS)
IN LINGUAGGIO HLS CHE POSSA ESSERE SINTETIZZATO:

- ABBIAMO USATO [HLS4ML](#) PER REALIZZARE
IMPLEMENTAZIONI FIRMWARE DI ALGORITMI ML



C/C++

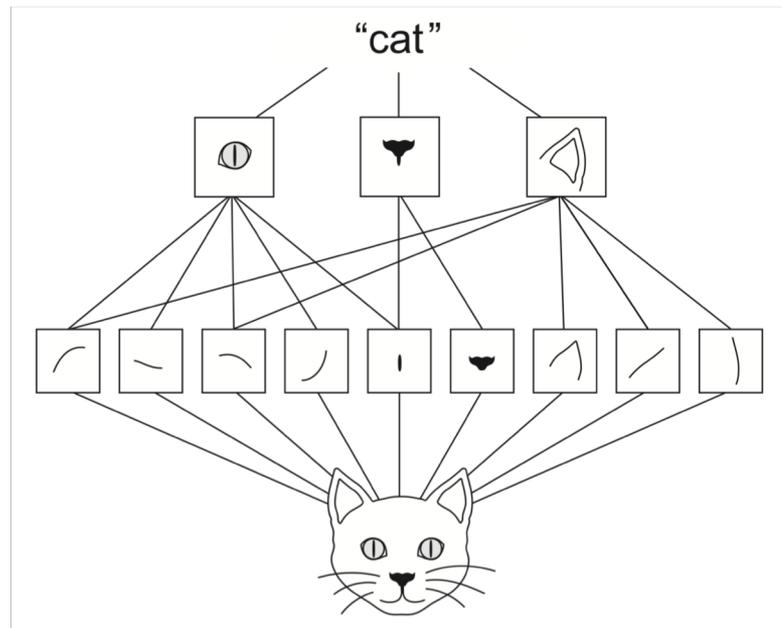
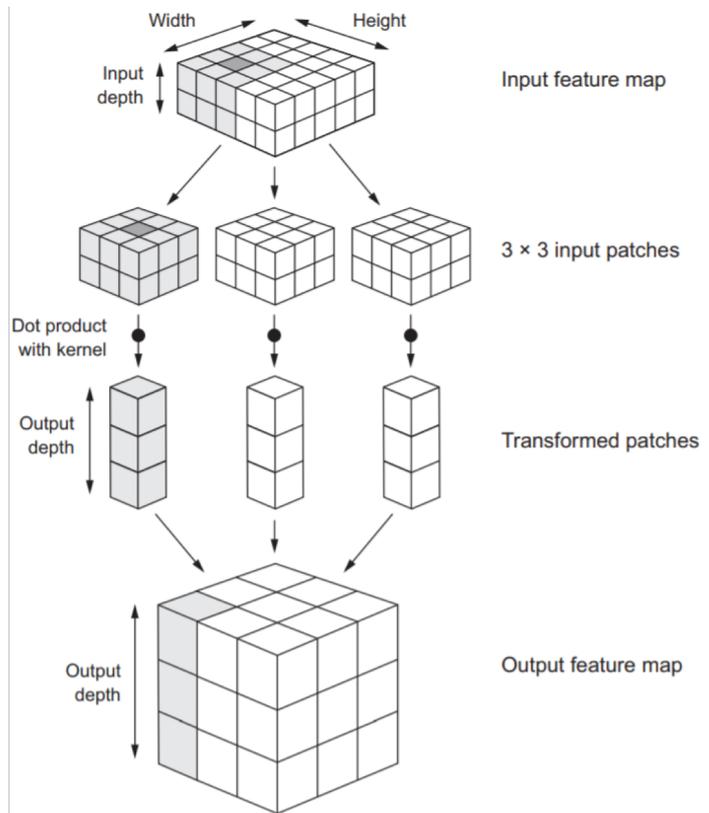
VIVADO HLS

VHDL/VERILOG

SINTESI ED IMPLEMENTAZIONE

LUNGHE OPERAZIONI
DI PROGRAMMAZIONE

CNN - FUNZIONAMENTO



ARCHITETTURA RETE DENSA

