

BIG DATA NEL PROTONE

STEFANO FORTE
UNIVERSITÀ DI MILANO & INFN



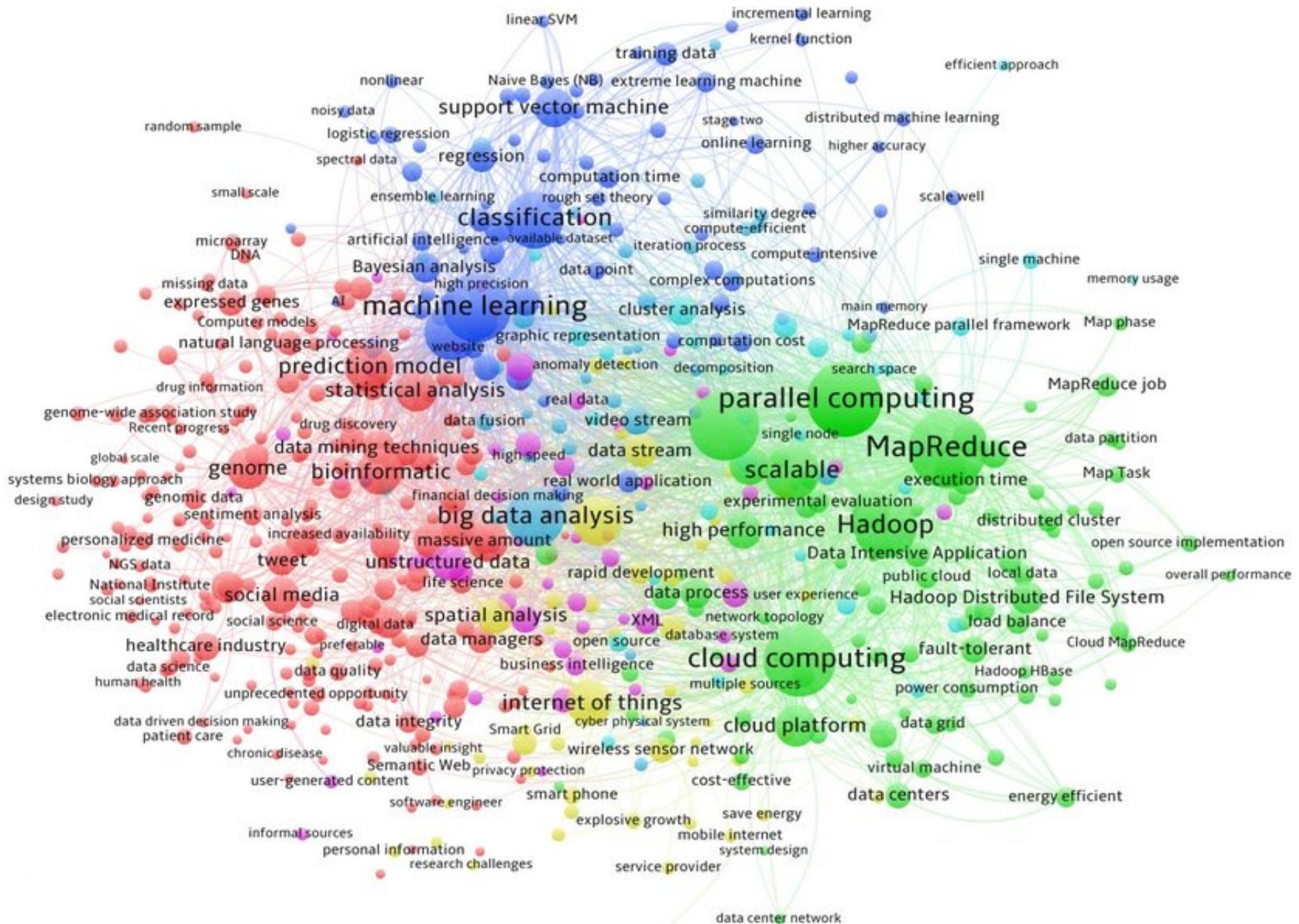
UNIVERSITÀ DEGLI STUDI DI MILANO
DIPARTIMENTO DI FISICA



PHYSICS DRINKS

MILANO, 27 APRILE 2022

BIG DATA?



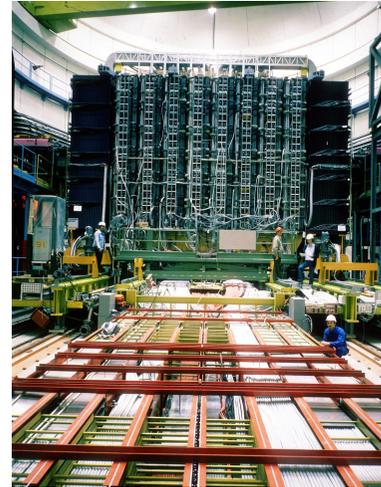
LA SCOPERTA DEI BOSONI W E Z (1984)

“BIG SCIENCE” MA NON ANCORA BIG DATA

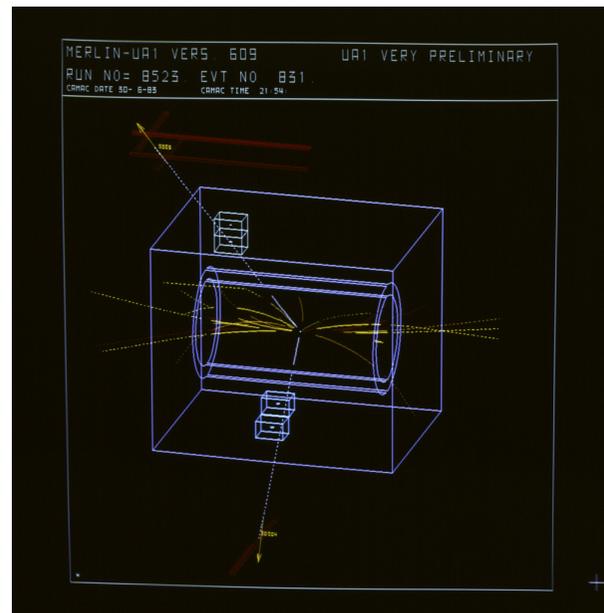
UA1: \sim 100 MEMBRI



UA1: lungh. 6m, diam. 2.3m



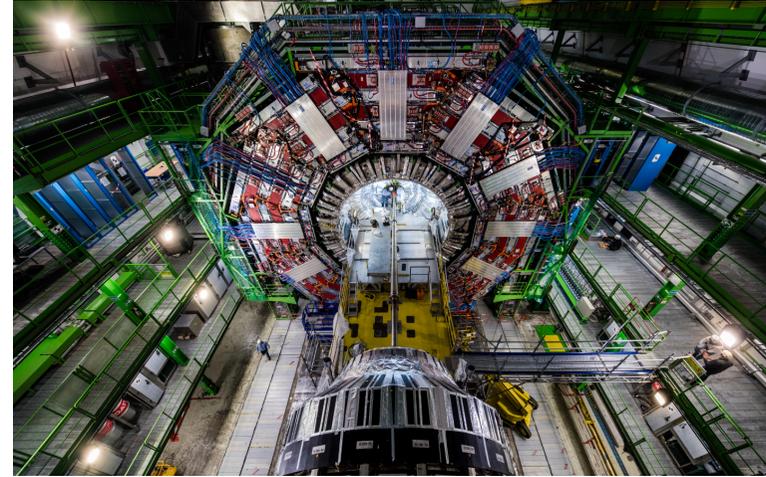
EVENTO: DECINE DI TRACCE



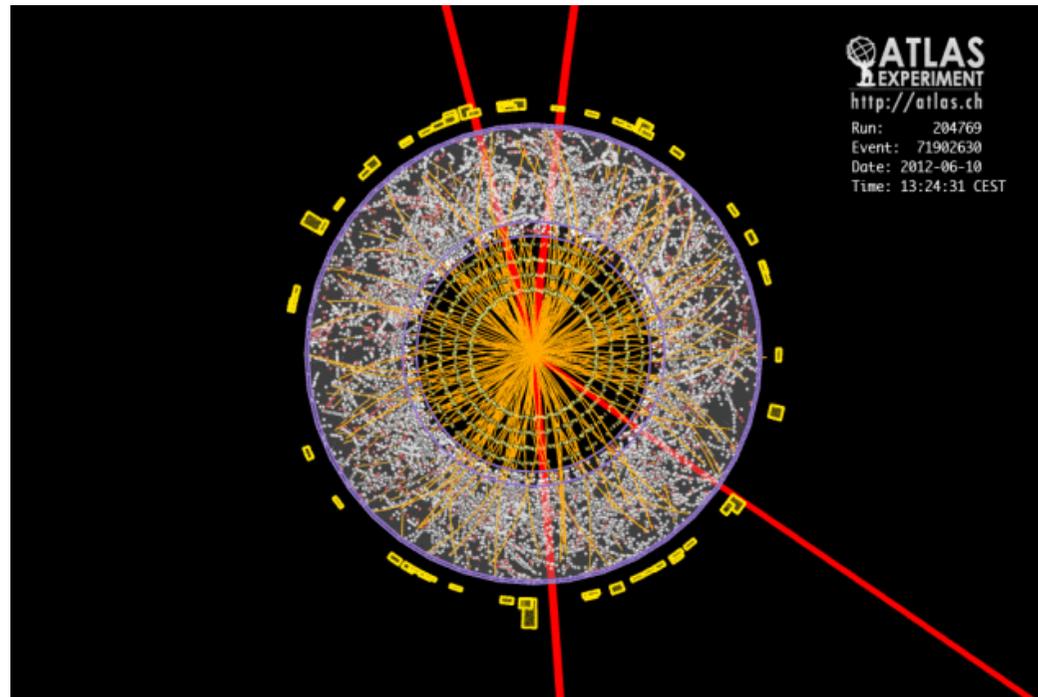
LA SCOPERTA DEL BOSONE DI HIGGS (2012)

CMS: ~ 2500 MEMBRI

CMS: lungh. 21m, diam. 15m

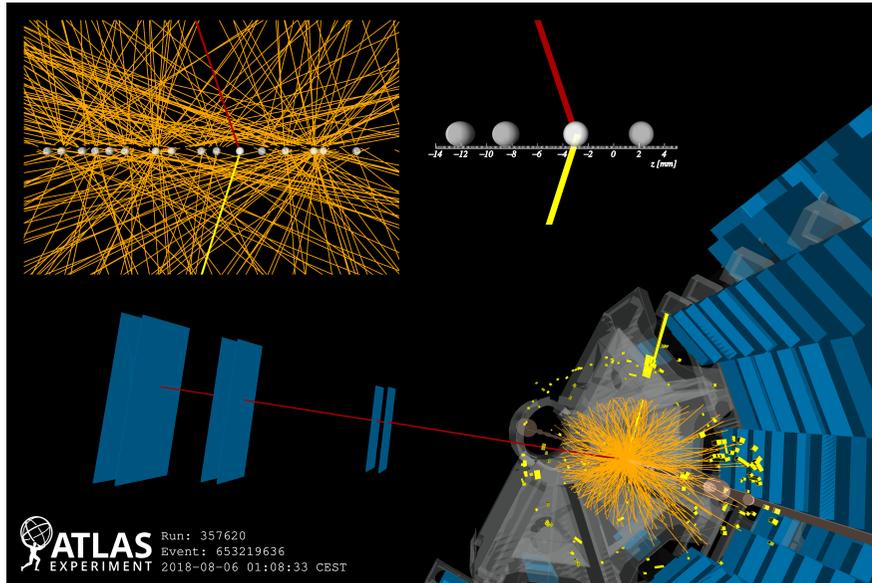


EVENTO: MIGLIAIA DI TRACCE

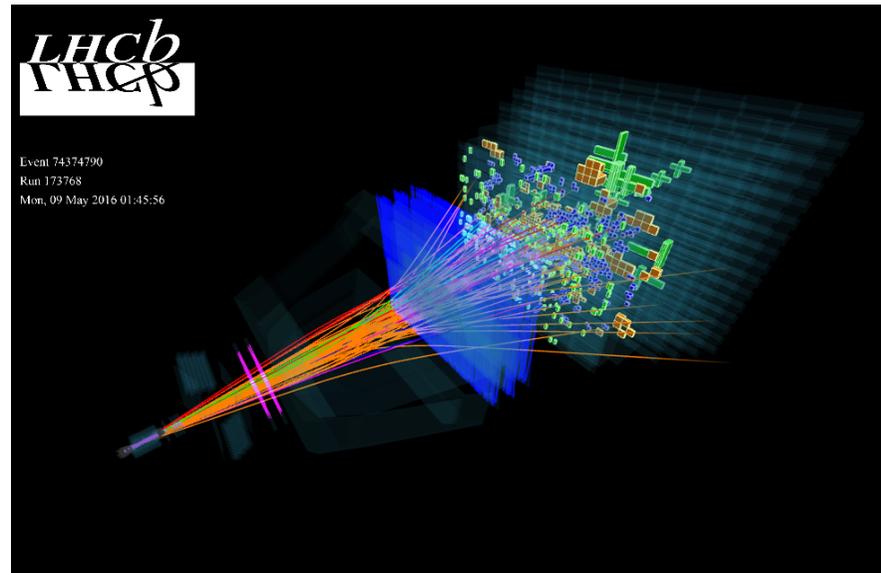


BIG DATA!

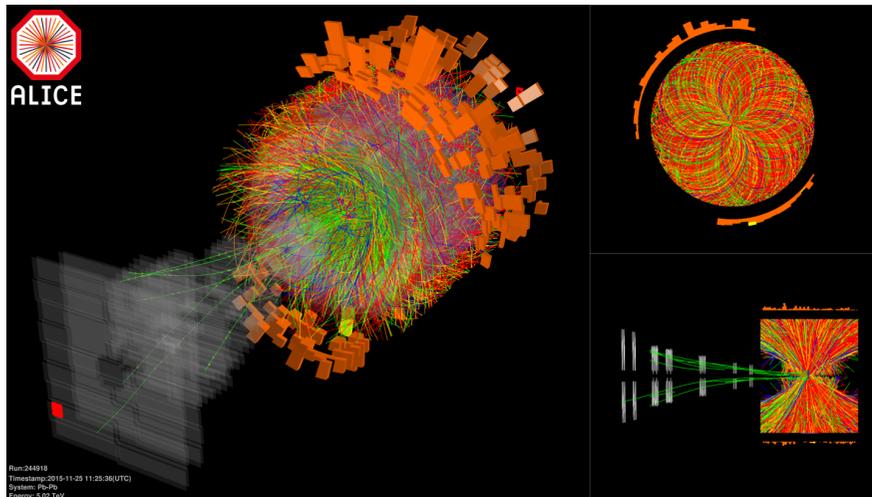
ATLAS



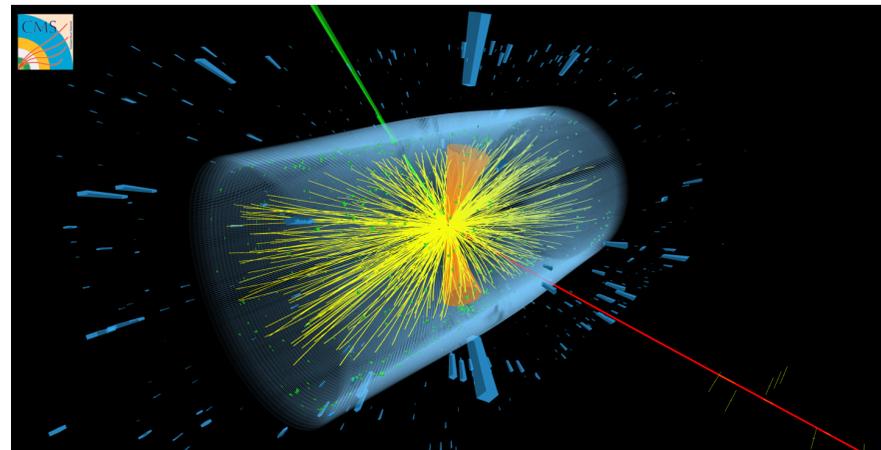
LHCb



ALICE

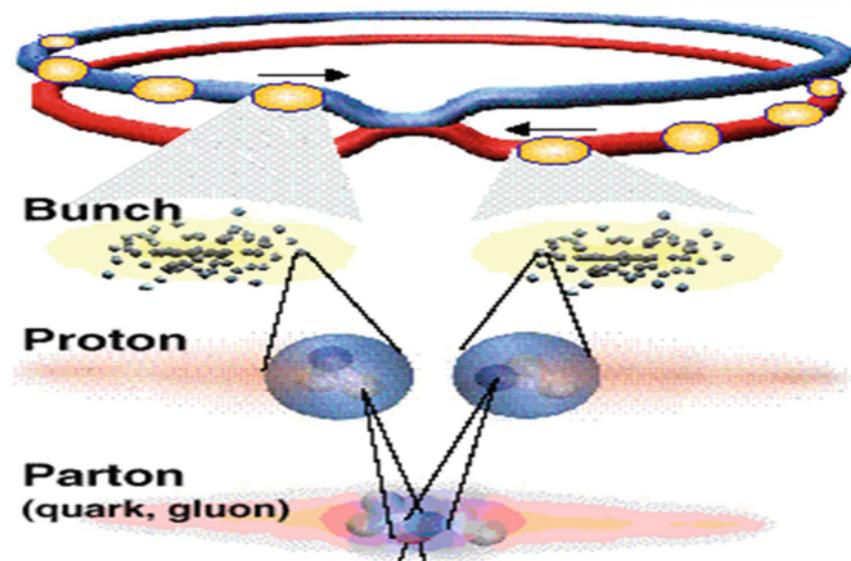
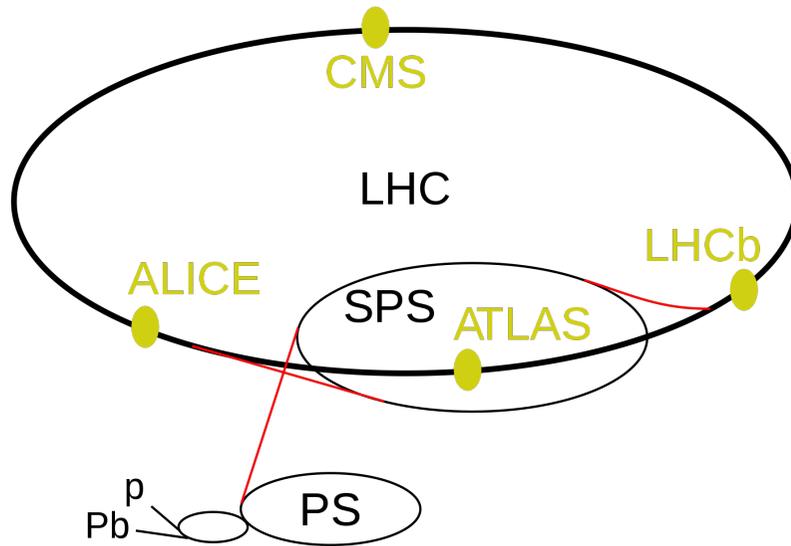


CMS



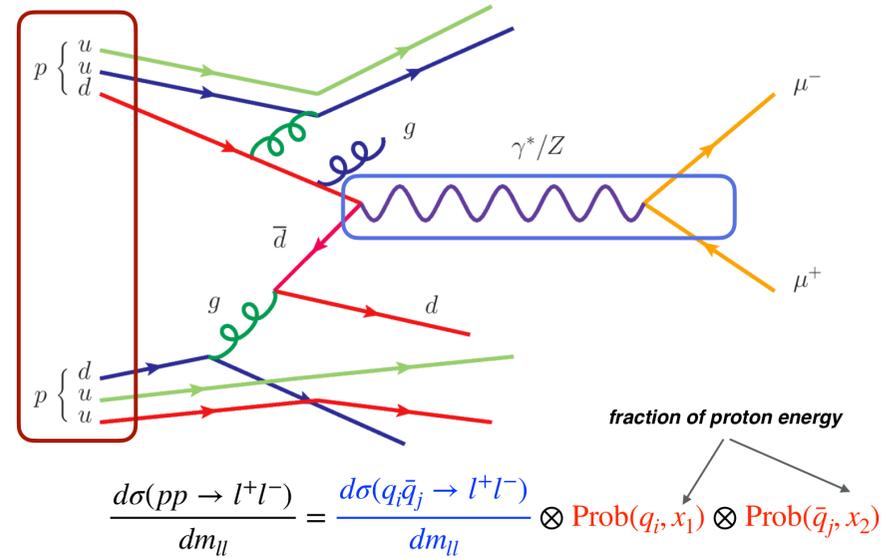
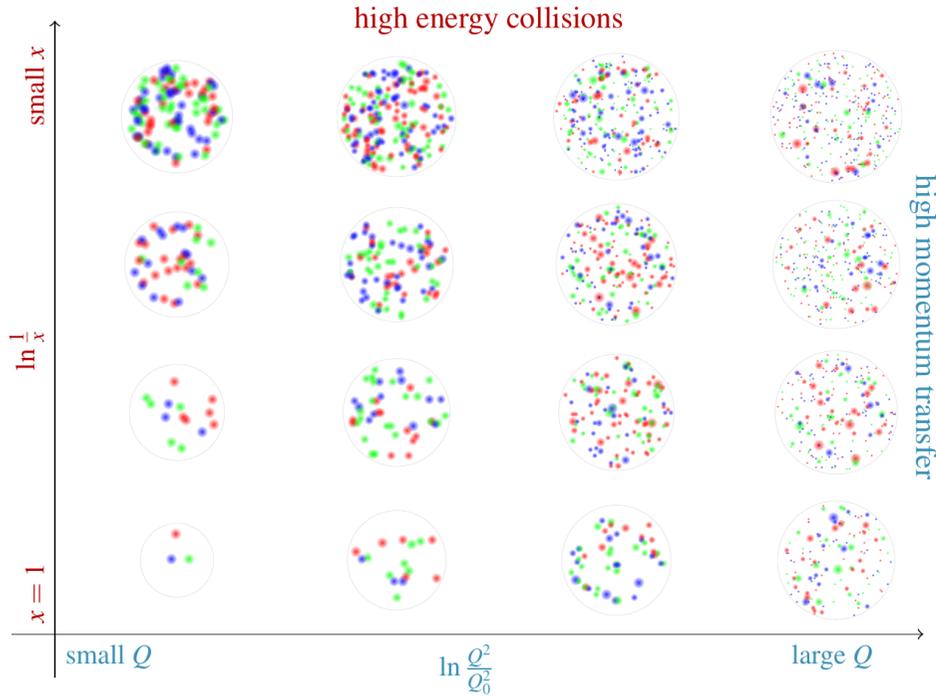
CIRCA 1 MILIARDO COLLISIONI/SEC; CIRCA 100 PETABYTE/ANNO

LHC: URTI FRA PROTONI



IL PROTONE NON È UNA PARTICELLA ELEMENTARE!

DENTRO IL PROTONE



PROBABILITÀ DI PROBABILITÀ:
 INFINITI² PARAMETRI
 DA ESTRARRE DAI DATI

⇔

QCD: UN SOLO
 PARAMETRO LIBERO Λ

SMALL DATA

IL PROTONE NEL 1984...



EUROPEAN ORGANIZATION FOR NUCLEAR RESEARCH

CERN-EP/85-108
11 July 1985

W PRODUCTION PROPERTIES AT THE CERN SPS COLLIDER

UA1 Collaboration, CERN, Geneva, Switzerland

Aachen¹ - Amsterdam (NIKHEF)² - Annecy (LAPP)³ - Birmingham⁴ - CERN⁵ -
Harvard⁶ - Helsinki⁷ - Kiel⁸ - London (Imperial College⁹ and Queen Mary College¹⁰) - Padua¹¹ -
Paris (Coll. de France)¹² - Riverside¹³ - Rome¹⁴ - Rutherford Appleton Lab.¹⁵ -
Saclay (CEN)¹⁶ - Victoria¹⁷ - Vienna¹⁸ - Wisconsin¹⁹ Collaboration

The corresponding experimental result for the 1984 data at $\sqrt{s} = 630$ GeV is

$$(\sigma \cdot B)_W = 0.63 \pm 0.05 (\pm 0.09) \text{ nb.}$$

This is in agreement with the theoretical expectation [14] of $0.47^{+0.14}_{-0.08}$ nb. We note that the 15%

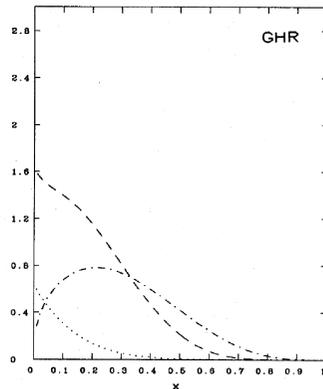


FIG. 25. Parton distributions of Glück, Hoffmann, and Reya (1982), at $Q^2=5$ GeV²; valence quark distribution $x[u_v(x)+d_v(x)]$ (dotted-dashed line), $xG(x)$ (dashed line), and $q_v(x)$ (dotted line).

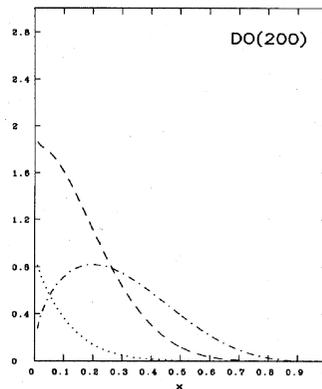


FIG. 27. "Soft-gluon" ($\Lambda=200$ MeV) parton distributions of Duke and Owens (1984) at $Q^2=5$ GeV²; valence quark distribution $x[u_v(x)+d_v(x)]$ (dotted-dashed line), $xG(x)$ (dashed line), and $q_v(x)$ (dotted line).

42

G. Altarelli et al. / Vector boson production

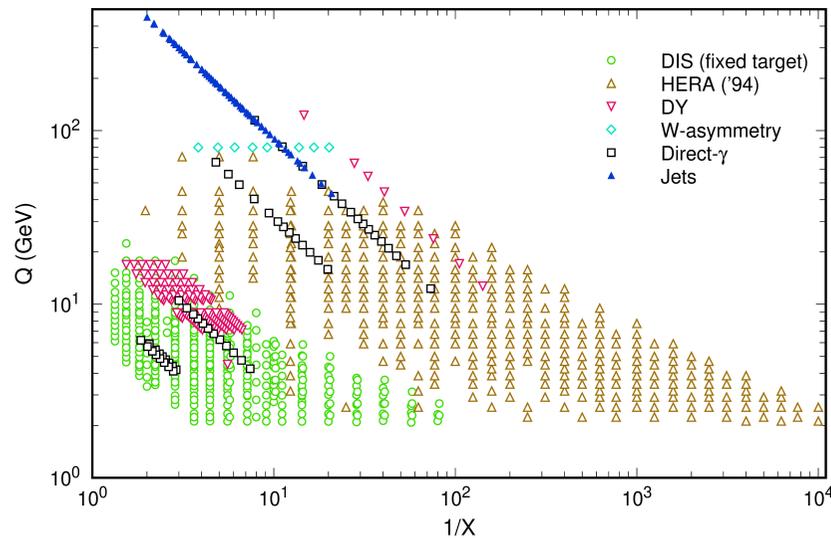
TABLE 2
Values (in nb) of the total cross sections for W^\pm and Z^0 production

\sqrt{s} (GeV)	$W^+ + W^-$		Z^0		$\frac{\sigma(W^+ + W^-)}{\sigma(Z^0)}$		$\frac{\sigma(W^+ + W^-)}{\sigma(Z^0)}$		$\frac{\sigma(W^+ + W^-)}{\sigma(Z^0)}$	
	GHR	DO1	DO2	GHR	DO1	DO2	GHR	DO1	DO2	
540	4.2	4.3	4.1	1.3	1.3	1.2	3.1	3.4	3.5	
700	6.2	6.3	6.1	2.0	1.9	1.8	3.1	3.3	3.4	
1000	9.5	9.5	9.6	3.1	3.0	2.9	3.1	3.2	3.3	
1300	12.5	12.5	12.9	4.0	3.9	3.9	3.1	3.2	3.3	
1600	15.5	15.6	16.5	5.0	4.8	5.0	3.1	3.2	3.3	

- SEMPLICE MODELLO CON DUE-TRE PARAMETRI
- NESSUNA STIMA DELL'ERRORE
- ACCURATEZZA QUALITATIVA

...E NEL 2000

DATASET CTEQ5 (1999)



- QUALCHE CENTINAIA DI DATI
- PARAMETRIZZAZIONE AD-HOC
- FIT MULTIPARAMETRICO
CON INCERTEZZA
- ACCURATEZZA SEMI-QUANTITATIVA

PARAMETRIZZAZIONE CTEQ5 (2006)

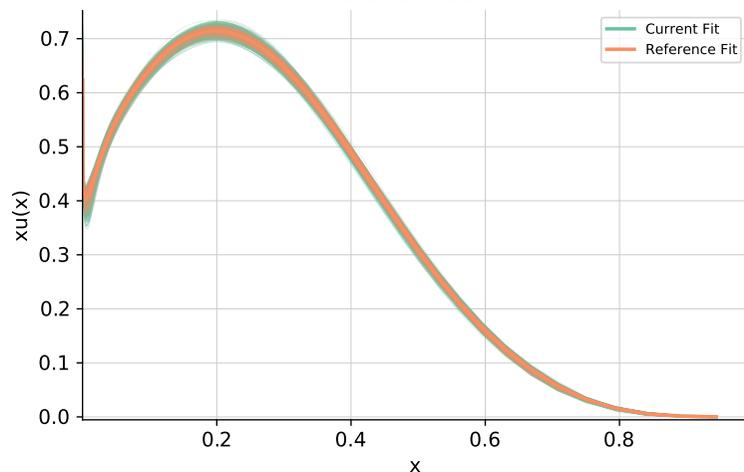
$$x f(x, Q_0) = A_0 x^{A_1} (1 - x)^{A_2} e^{A_3 x} (1 + e^{A_4} x)^{A_5}, \quad (6 \text{ funzioni}, 22 \text{ parametri})$$

TOWARDS BIG DATA

NNPDF (2002-2017)

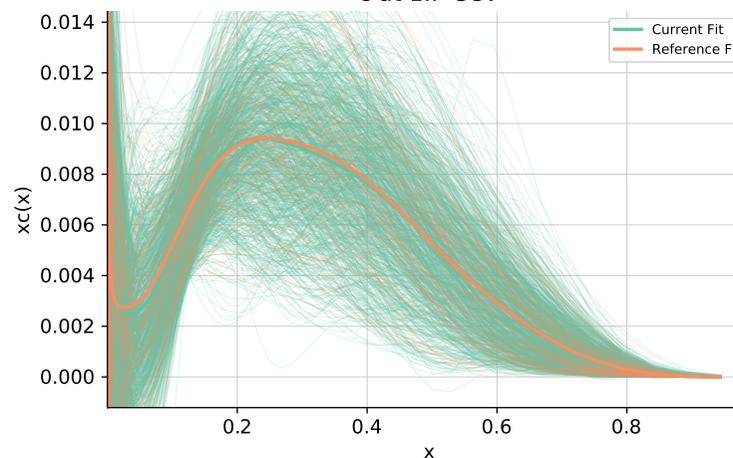
QUARK UP

u at 1.7 GeV



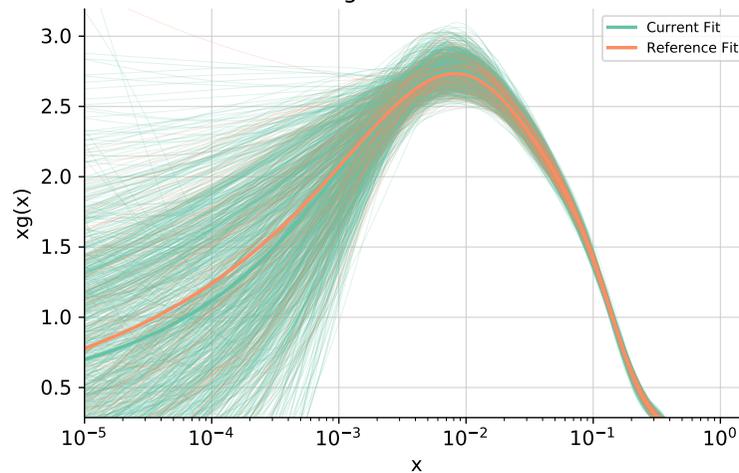
QUARK CHARM

c at 1.7 GeV



GLUONE

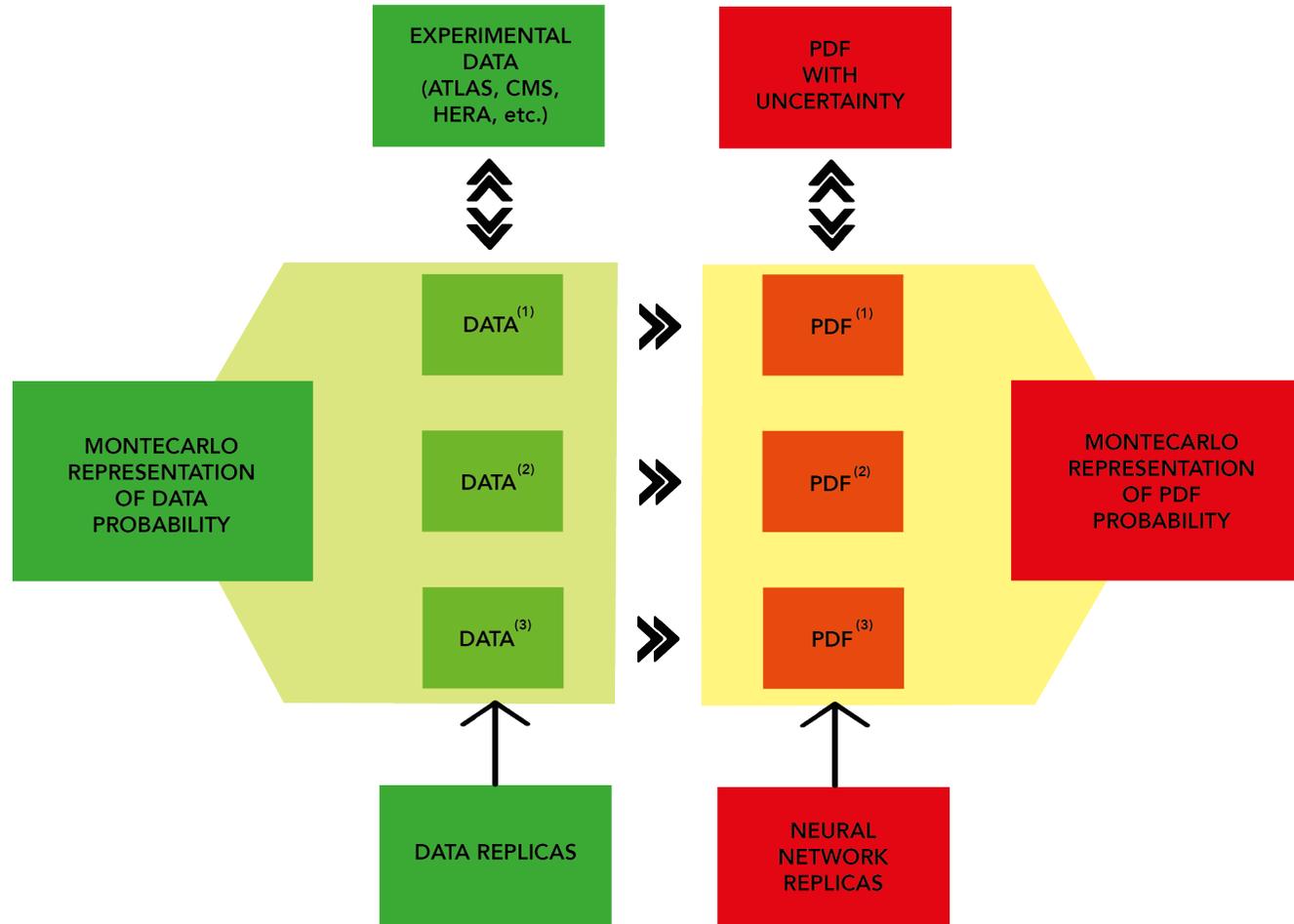
g at 1.7 GeV



RAPPRESENTAZIONE MONTECARLO \Leftrightarrow DISTRIBUZIONE DI PROBABILITÀ

IL MONTECARLO FUNZIONALE

REPLICHE DI FUNZIONI \Leftrightarrow PROBABILITÀ IN UNO SPAZIO DI FUNZIONI

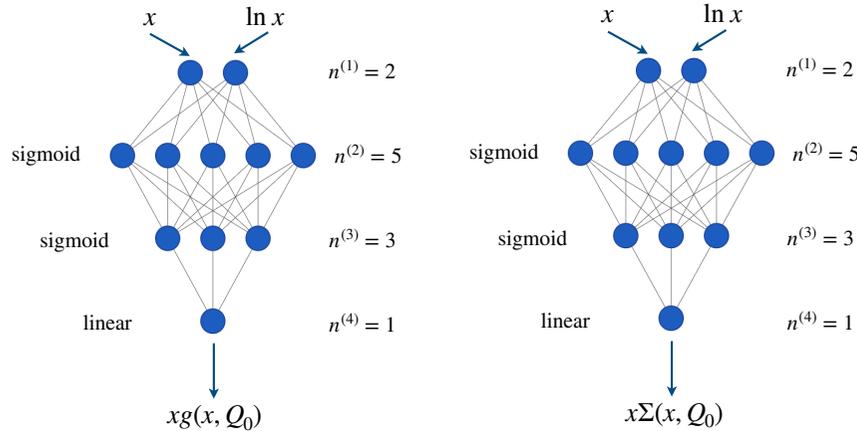


PDF OUTPUT: $f_i^{(a)}(x, \mu)$; $i = \text{up, antiup, down, antidown, strange, antistrange, charm, gluone}$;

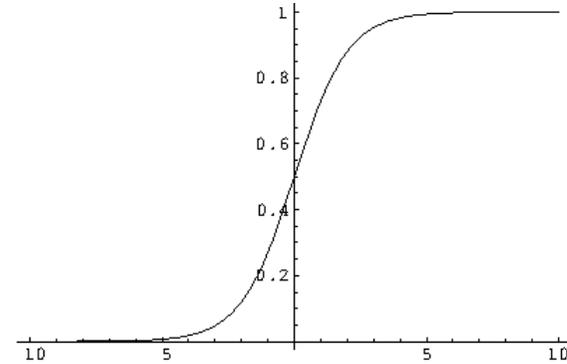
$j = 1, 2, \dots, N_{\text{rep}}$

RETI NEURALI

ARCHITETTURA



FUNZIONE DI ATTIVAZIONE



- INTERPOLANTE UNIVERSALE

- PUÒ RIPRODURRE

QUALUNQUE FORMA

FUNZIONALE

- COMPLESSITÀ CRESCE

DURANTE L'ADDESTRAMENTO

$$F_{\text{out}}^{(i)}(\vec{x}_{\text{in}}) = F\left(\sum_j \omega_{ij} x_{\text{in}}^j - \theta_i\right)$$

PARAMETERI

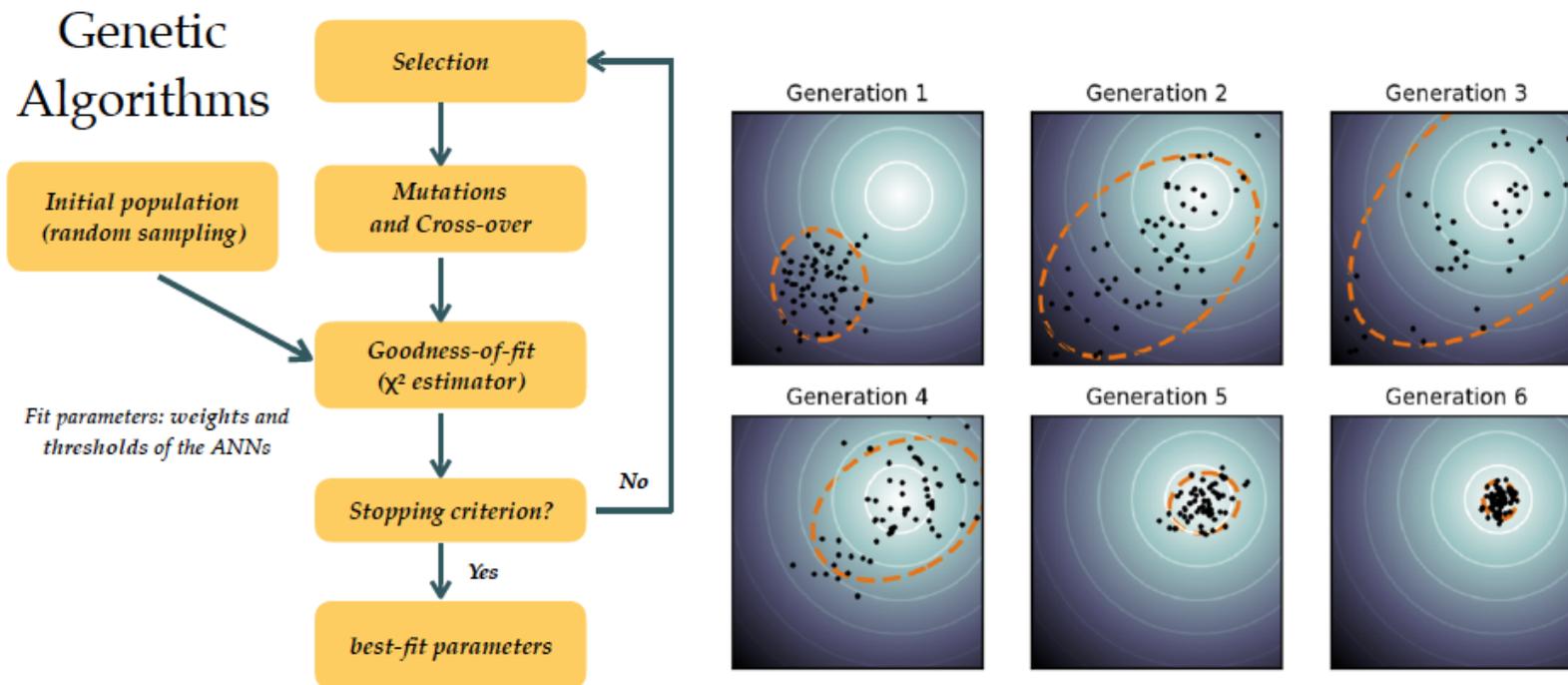
- PESI ω_{ij}

- SOGLIE θ_i

NNPDF1.0-3.1: 2 - 5 - 3 - 1 NN PER OGNI PDF: $37 \times 8 = 296$ PARAMETETRI

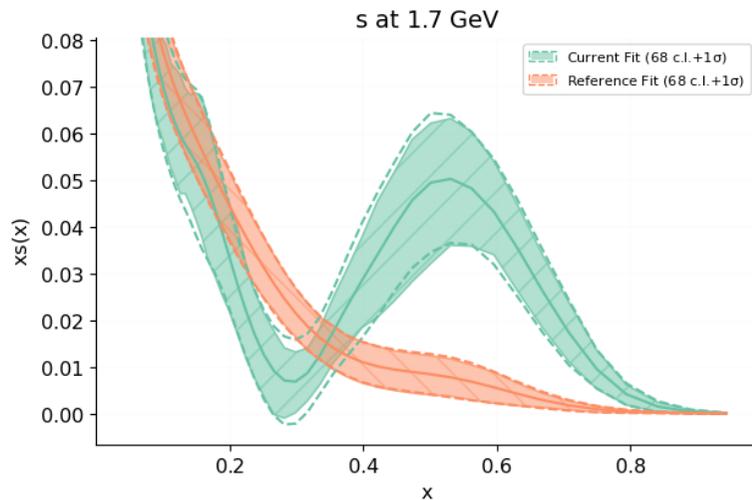
ALGORITMI GENETICI

- MUTAZIONE CASUALE DEI PARAMETRI DELLA RETE
- SELEZIONE DEL PIÙ ADATTO

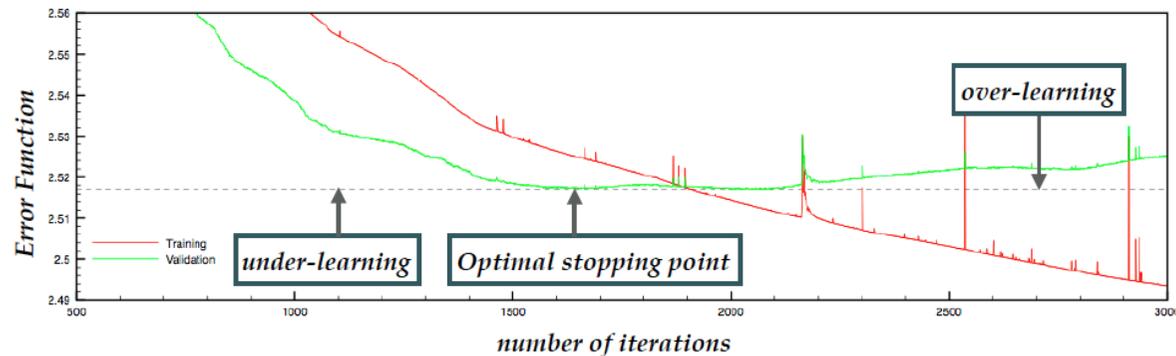


CONVALIDA INCROCIATA

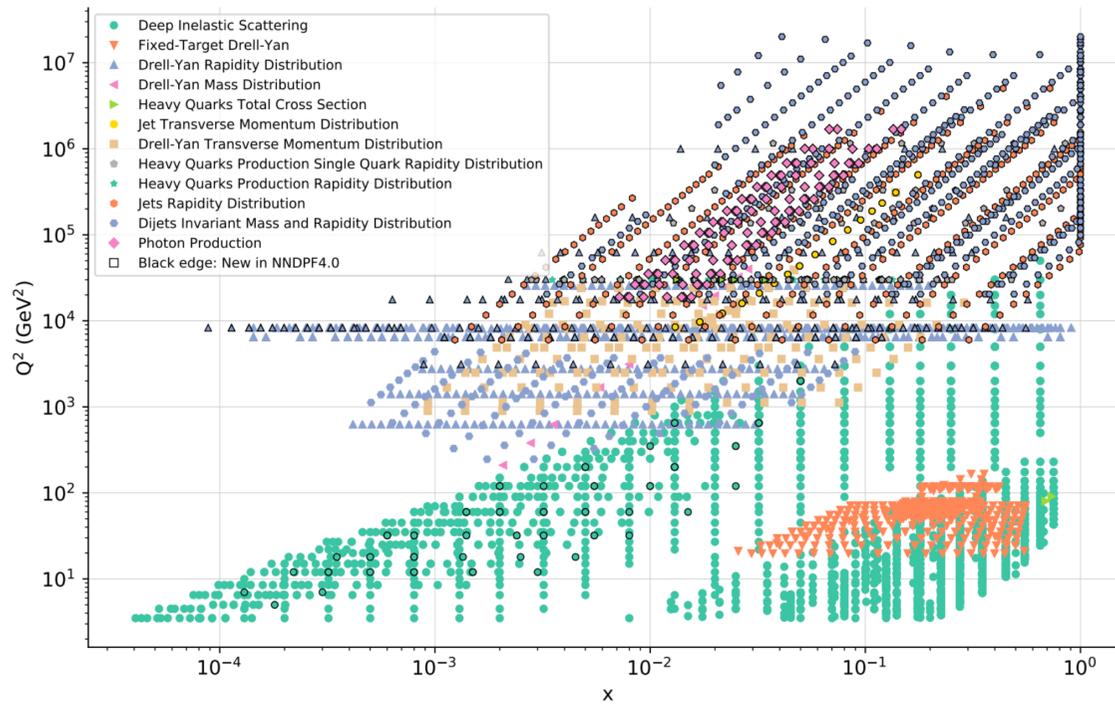
- LA RETE NEURALE
PUÒ RIPRODURRE TUTTO
- ANCHE IL RUMORE!



- DATI DIVISI CASUALMENTE
DUE GRUPPI: ADDESTRAMENTO
E CONVALIDA
- RETI NEURALI OTTIMIZZATE SUI
DATI DI ADDESTRAMENTO
- QUALITÀ DETERMINATA DAL
GRUPPO DI CONVALIDA



DATI

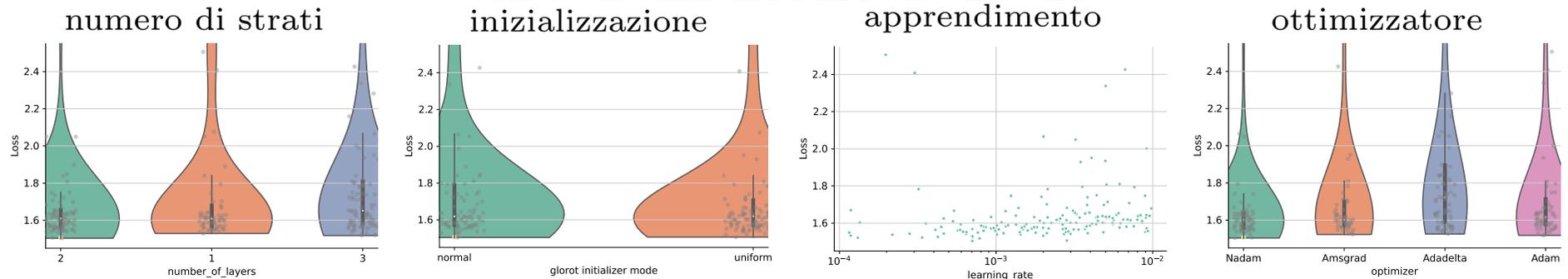


- ~ 50 ESPERIMENTI: ~ 5000 DATI SPERIMENTALI
- 1000 REPLICHE, 8 RETI NEURALI
- ~ 20000 GENERAZIONI ALGORITMO GENETICO
- 80 MUTANTI PER GENERAZIONE
- $\sim 10^{13}$ PREDIZIONI TEORICHE PER UNA DETERMINAZIONE TIPICA

BIG DATA!

IPEROTTIMIZZAZIONE

CHI SCEGLIE LA METODOLOGIA?
inizializzazione apprendimento



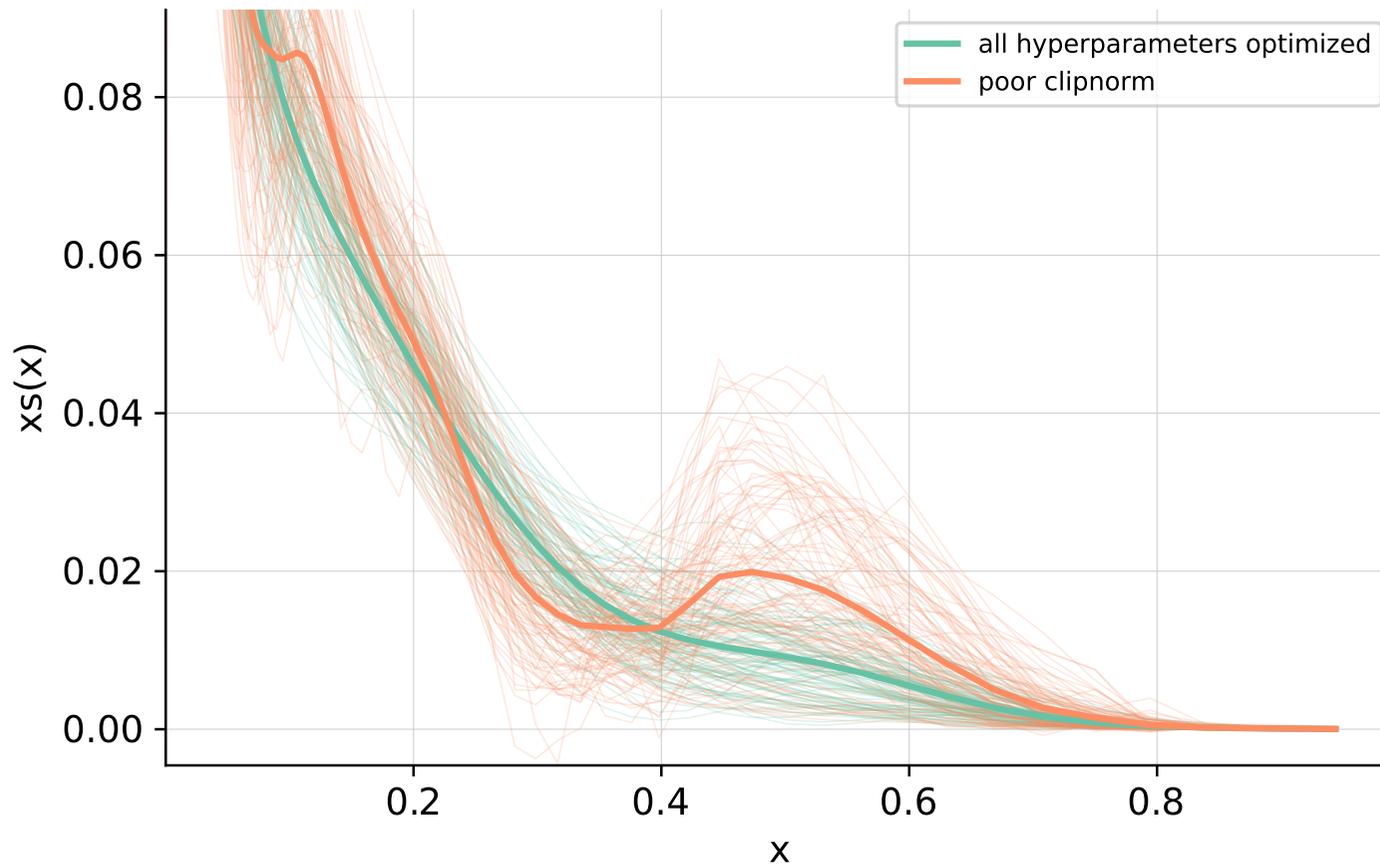
HYPEROPT PARAMETERS

NEURAL NETWORK	FIT OPTIONS
NUMBER OF LAYERS	OPTIMIZER
SIZE OF EACH LAYER	INITIAL LEARNING RATE
DROPOUT	MAXIMUM NUMBER OF EPOCHS
ACTIVATION FUNCTIONS	STOPPING PATIENCE
INITIALIZATION FUNCTIONS	POSITIVITY& INTEGRABILITY MULTIPLIER

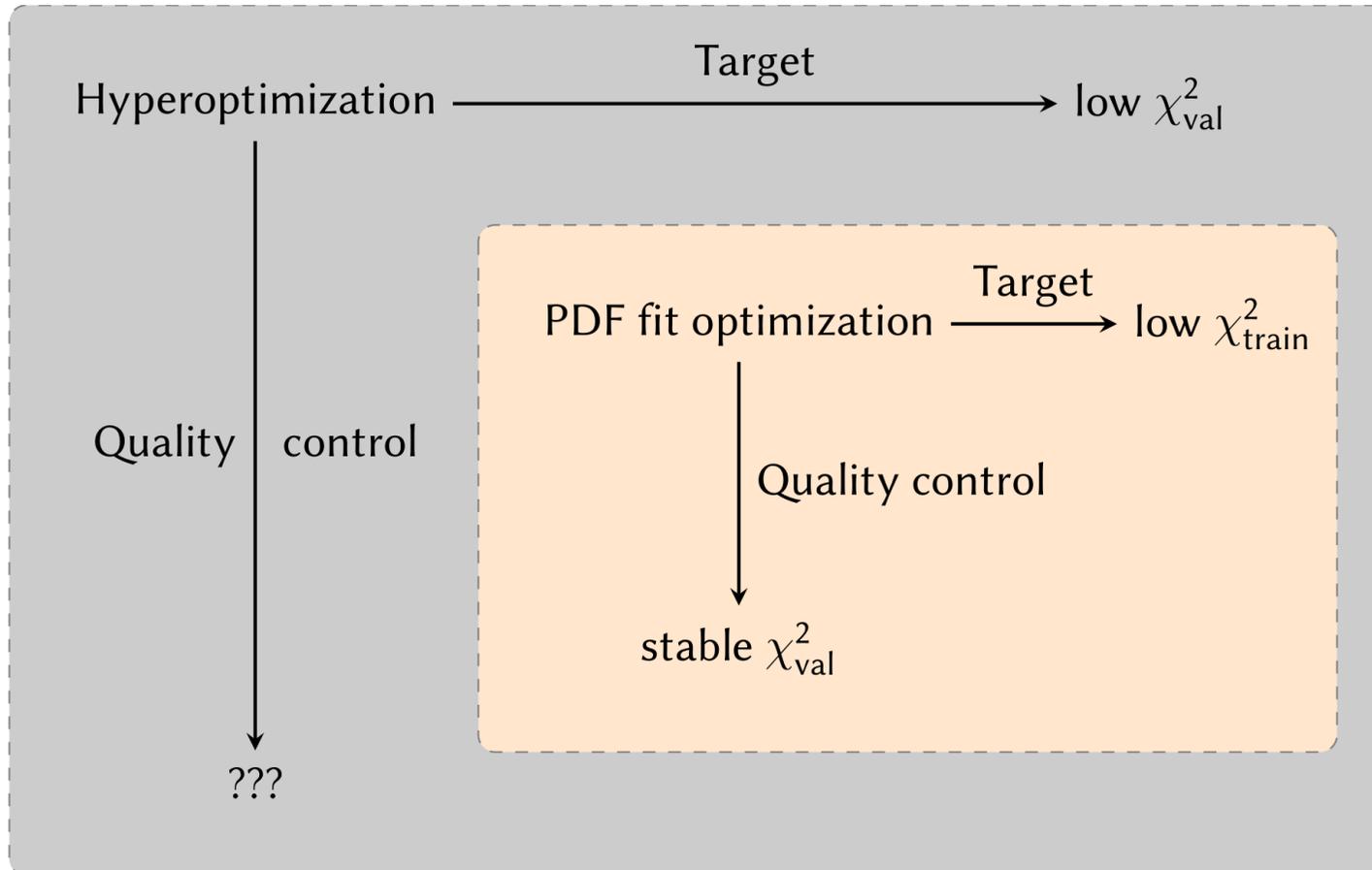
OVERFITTING

METODOLOGIA TROPPO AGGRESSIVA

s at 1.7 GeV

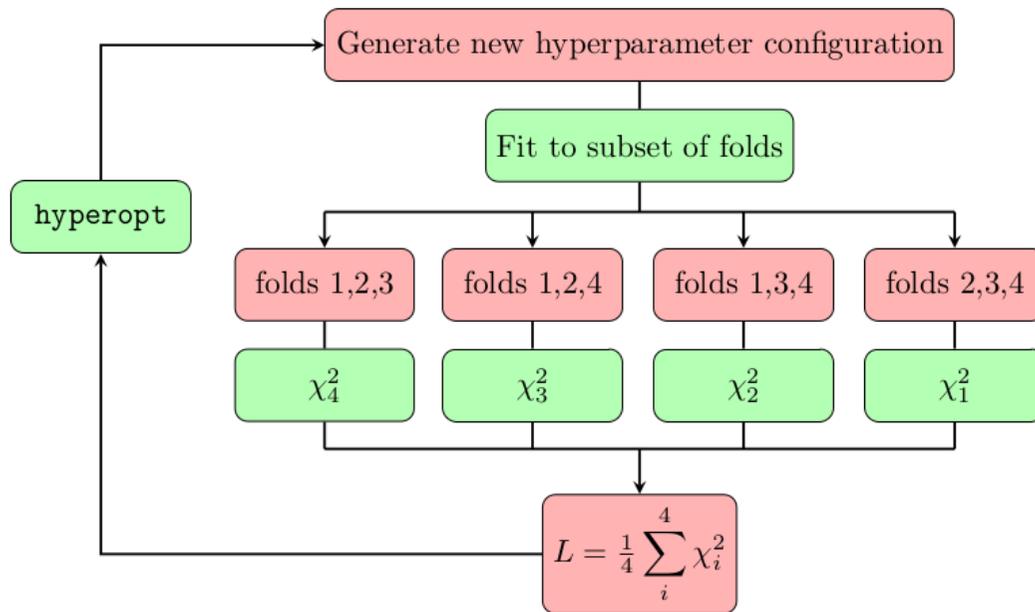


QUAL È LA METODOLOGIA OTTIMALE?



MANCA IL CRITERIO DI QUALITÀ

K-FOLDING



- DATI **SUDDIVISI**
IN n GRUPPI (FOLDS)
- UN GRUPPO **ESCLUSO**
A TURNO
- FIT MIGLIORE
⇒ **GENERALIZZA**
CORRETTAMENTE AI DATI
ESCLUSI

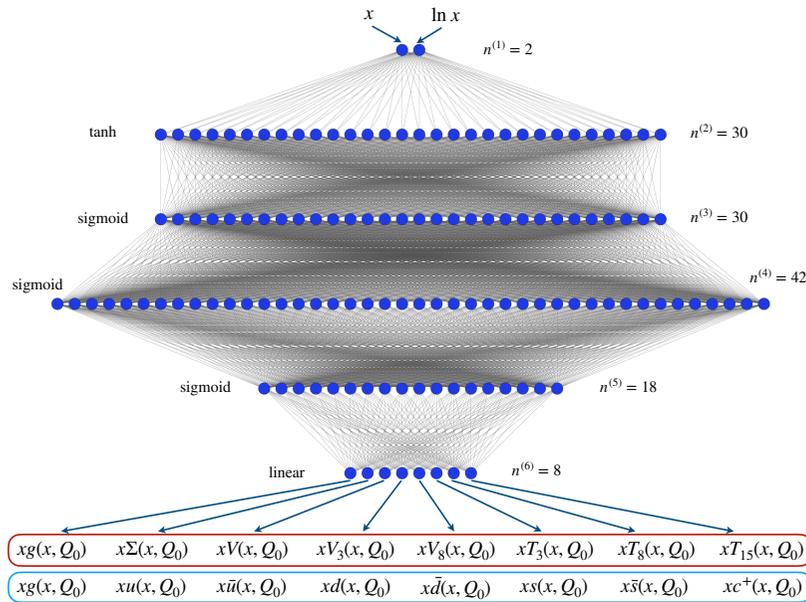
RISULTATI

15000 CICLI DI IPEROTTIMIZZAZIONE; 4 FOLDS

FINAL NNPDF4.0 SET-UP

AFTER HYPEROPTIMIZATION

UN'ARCHITETTURA TESTATA



PARAMETER	VALUE
ARCHITECTURE	25-20-8
ACTIVATION	HYPERBOLIC TANGENT
INITIALIZER	glorot_normal
OPTIMIZER	Nadam
CLIPNORM	$6 \cdot 10^{-6}$
LEARNING RATE	$2.6 \cdot 10^{-3}$
MAX EPOCHS	$17 \cdot 10^3$
STOPPING PATIENCE	10% OF MAX EPOCHS
INITIAL POSITIVITY λ	185
INITIAL INTEGRABILITY λ	10

IL PROTONE 2022

IL CODICE

AVERAGE FITTING TIME PER REPLICAS AND USE OF RESOURCES

SAME DATASET FOR OLD AND NEW METHODOLOGIES IN CPU AND GPU

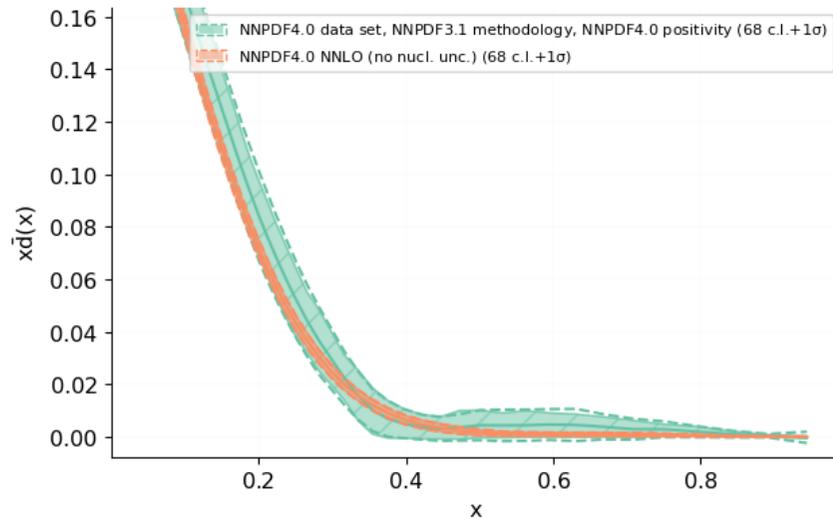
CPU: INTEL(R) CORE(TM) I7-4770 AT 3.40GHZ; GPU: NVIDIA TITAN V

	NNPDF31 CODEBASE	NNPDF40 CODEBASE IN CPU	NNPDF40 CODEBASE IN GPU
TIME	15.2 H.	38 ± 5 MIN.	6.6 MIN.
RAM USE	1.5 GB	6.1 GB	NA

LE PDF

ANTIDOWN: NNPDF3.1 VS NNPDF4.0

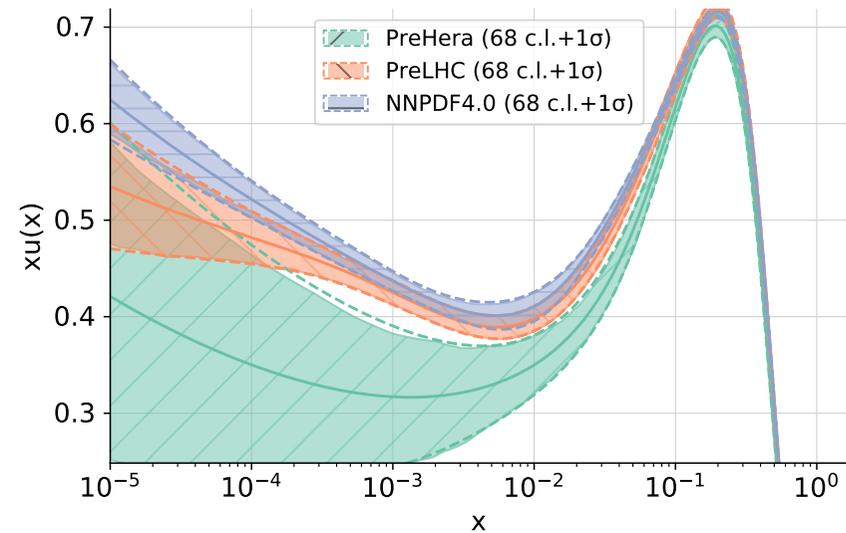
\bar{d} at 1.7 GeV



IL "FUTURE TEST"

QUARK UP: DATI 1995, 2005, 2021

u at 1.7 GeV



NEWS

IMPATTO



NNPDF_Team @NNPDF_Team · Apr 25



A nice way to celebrate the 20th anniversary of the paper that started it all: [hep-ph/0204232](#). [#H2020](#) [#EUfunded](#)



INSPIRE HEP @inspirehep · Apr 25

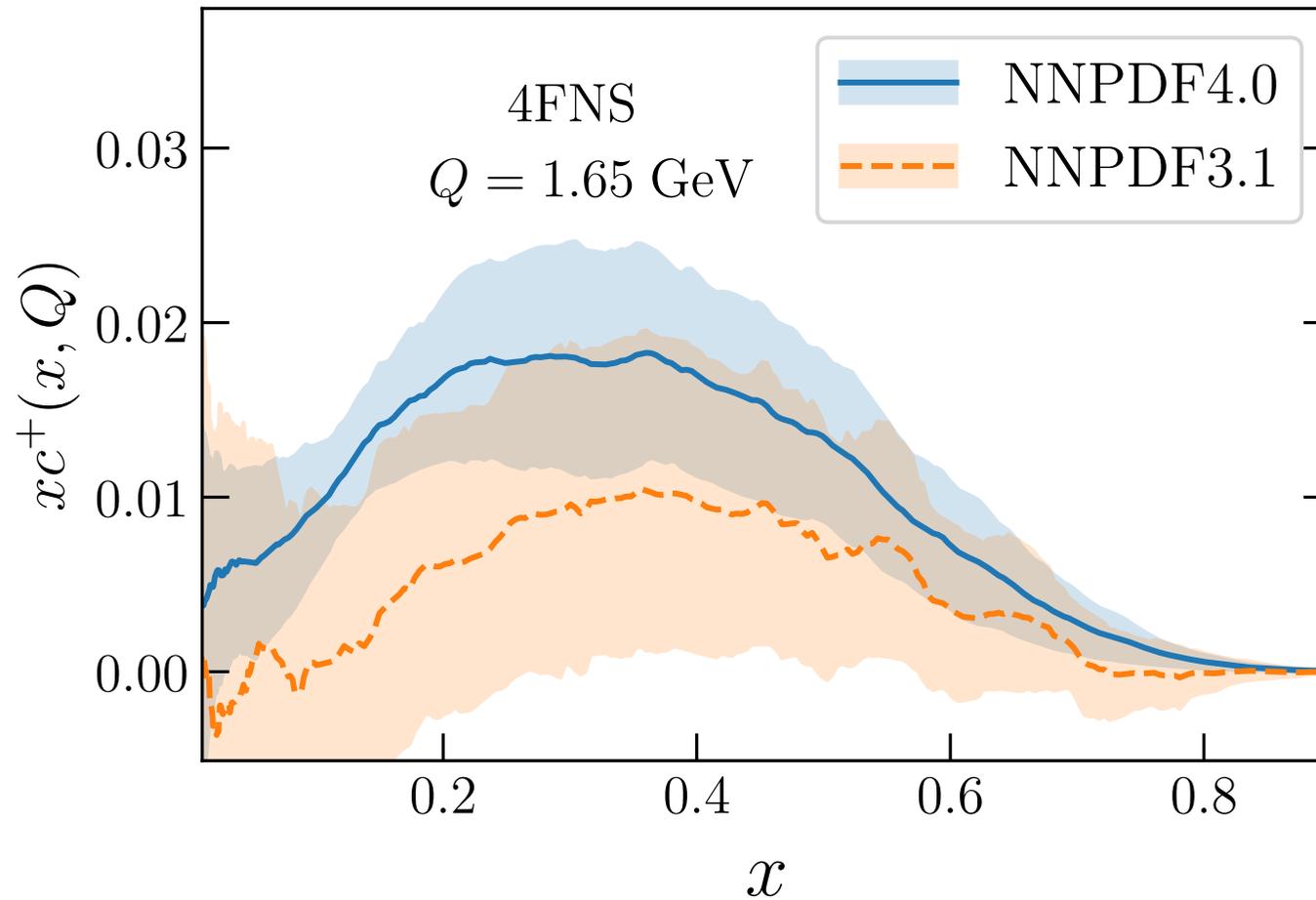
.@NNPDF_Team's 2015 JHEP article
"Parton distributions for the LHC Run II"
[inspirehep.net/literature/132...](#)
reaches 3,000 citations.

[#topcites](#) @SpringerPhysics



IL “CHARM” DEL PROTONE

VECCHIA VS. NUOVA METODOLOGIA



BIGGER DATA?

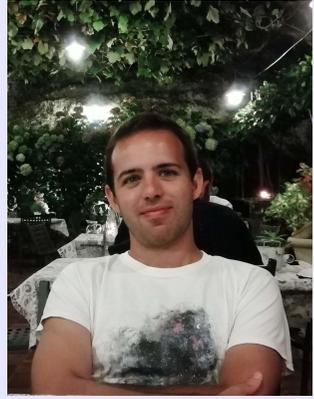
WHAT NEXT?

- REINFORCEMENT LEARNING
- IMPARARE LA DISTRIBUZIONE DI PROBABILITÀ
- ELEVATA PARALLELIZZAZIONE

A. Barontini



A. Candido



S. Carrazza



J. Cruz Martinez



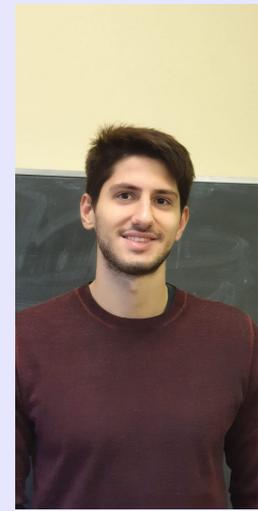
F. Hekhorn



K. Kudashkin



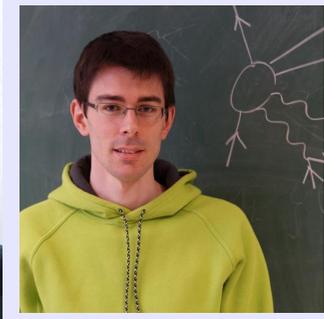
N. Laurenti



T. Rabemananjara



C. Schwan



R. Stegeman

