

European Research Council
Established by the European Commission



BIG DATA

STEFANO FORTE
UNIVERSITÀ DI MILANO & INFN



UNIVERSITÀ DEGLI STUDI DI MILANO
DIPARTIMENTO DI FISICA

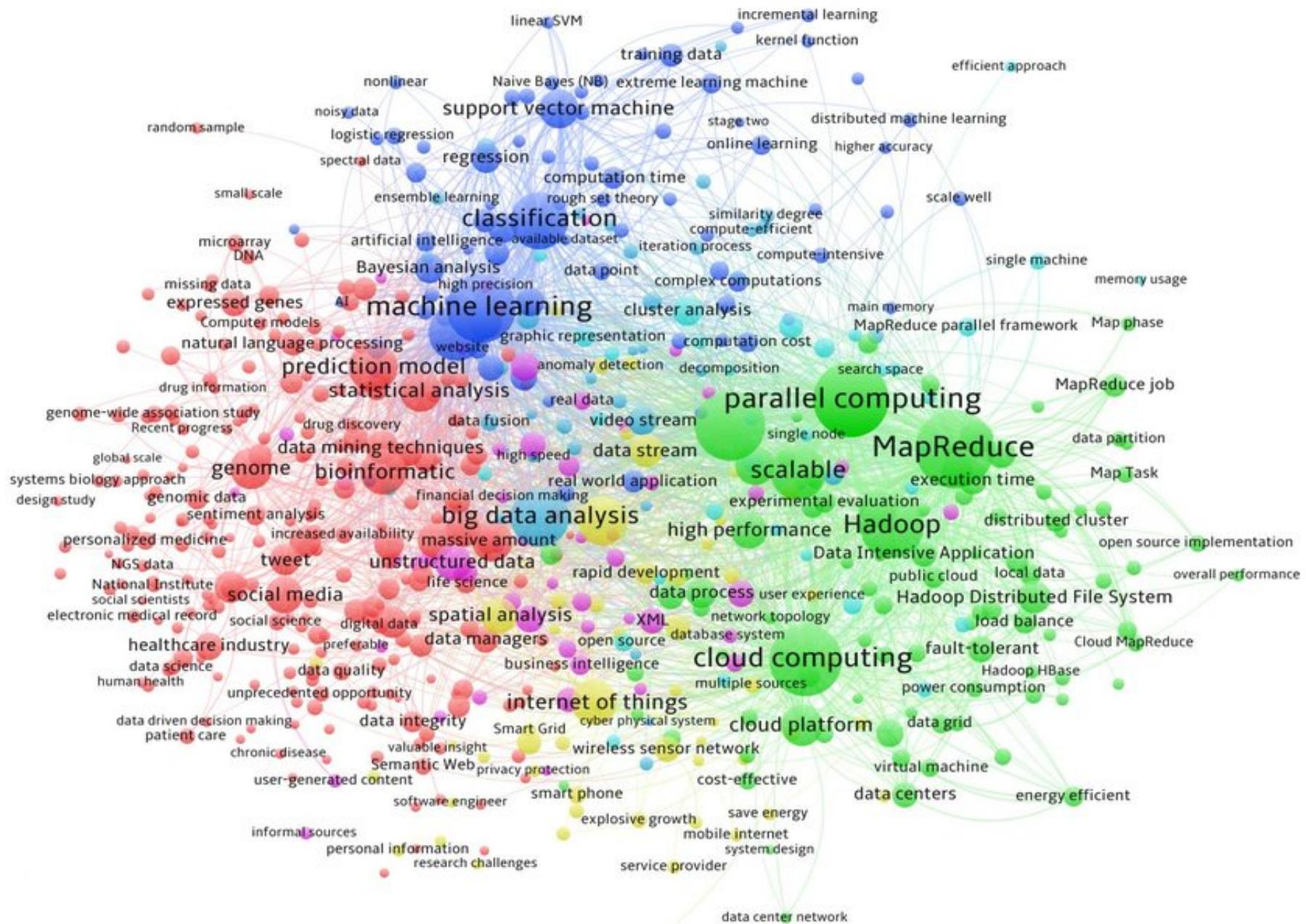


SULLE SPALLE DEI GIGANTI

BOLOGNA, 20 MAGGIO 2021

This project has received funding from the European Union's Horizon 2020 research and innovation programme under grant agreement No 740006

BIG DATA?



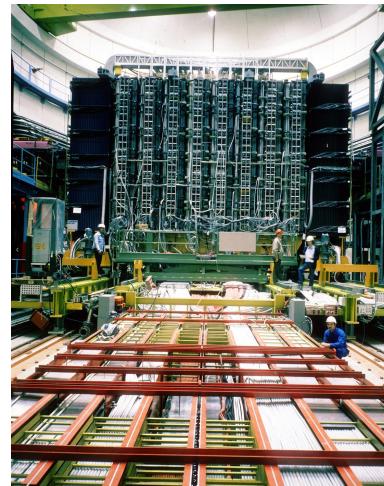
LA SCOPERTA DEI BOSONI W E Z (1984)

“BIG SCIENCE” MA NON ANCORA BIG DATA

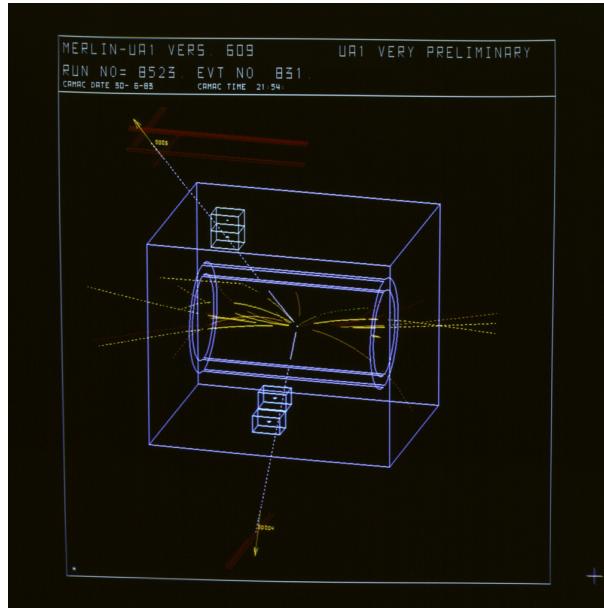
UA1: ~ 100 MEMBRI



UA1: lungh. 6m, diam. 2.3m



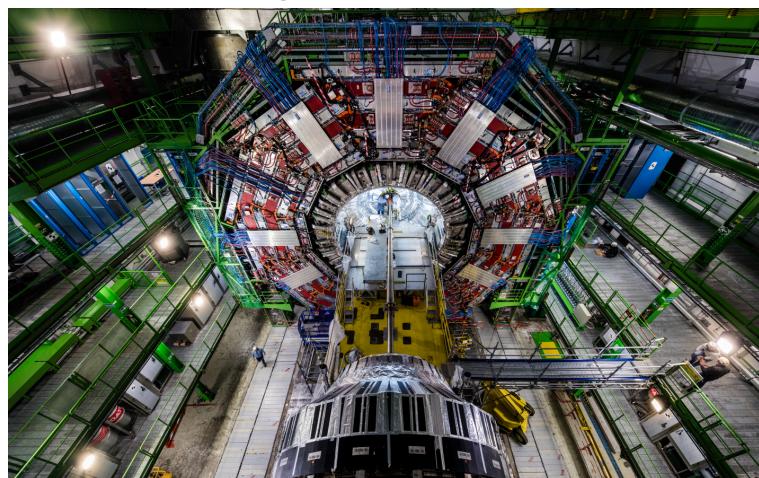
EVENTO: DECINE DI TRACCE



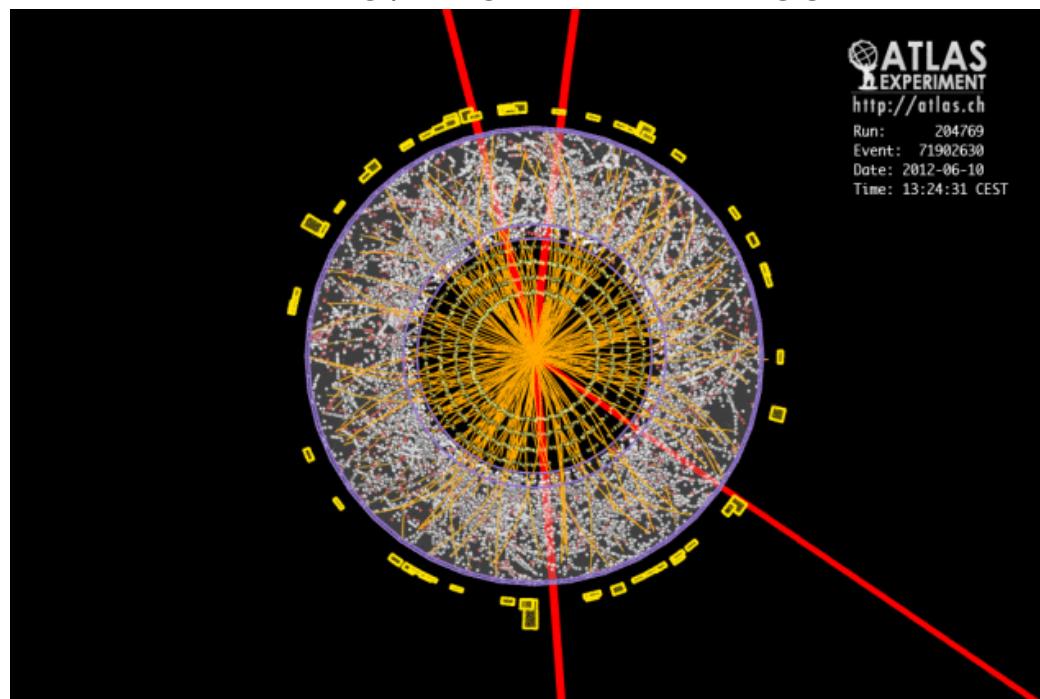
LA SCOPERTA DEL BOSONE DI HIGGS (2012)

CMS: \sim 2500 MEMBRI

CMS: lungh. 21m, diam. 15m

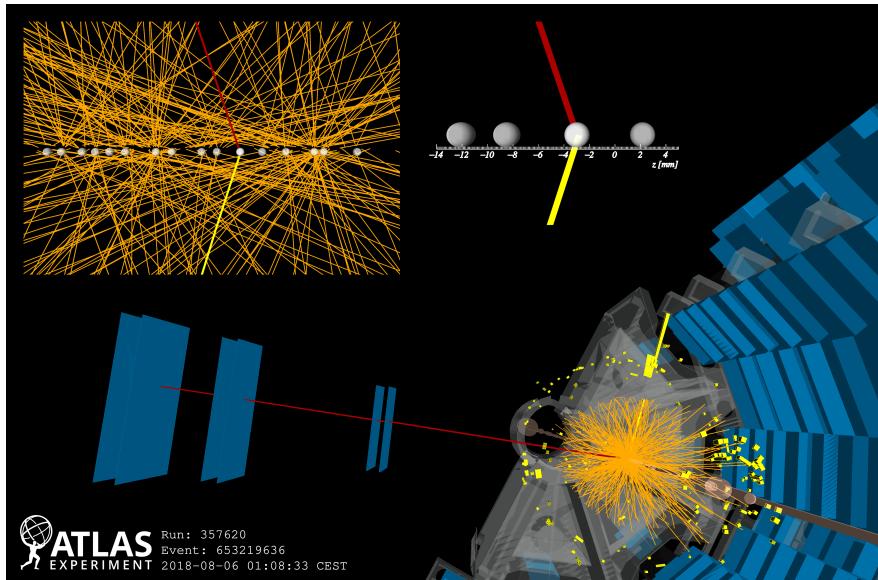


EVENTO: MIGLIAIA DI TRACCE

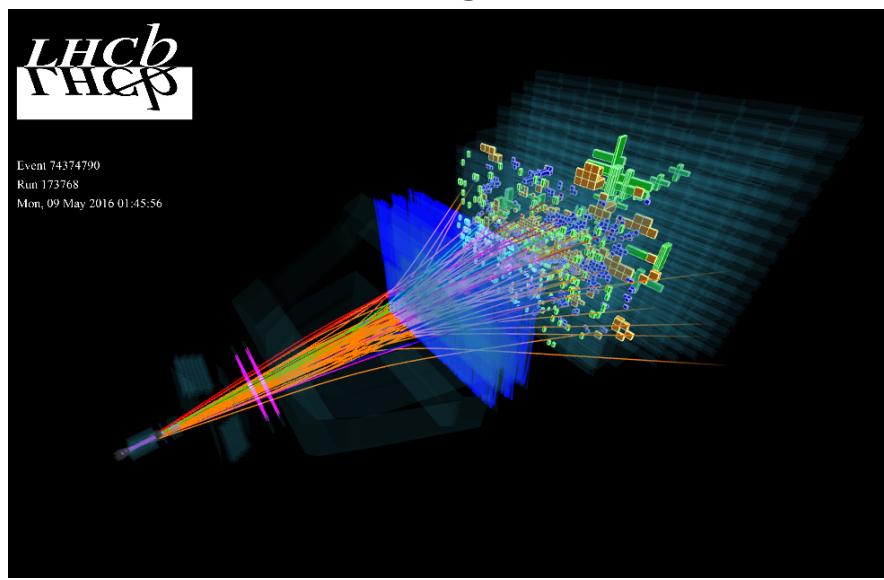


BIG DATA!

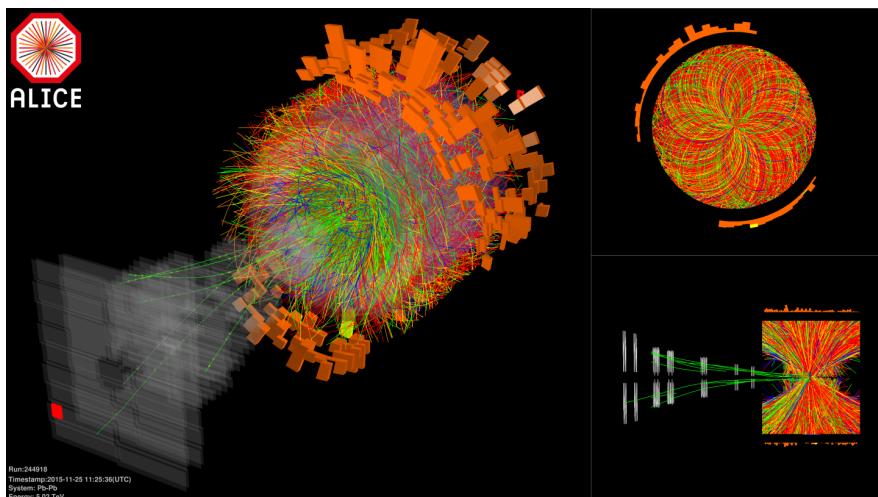
ATLAS



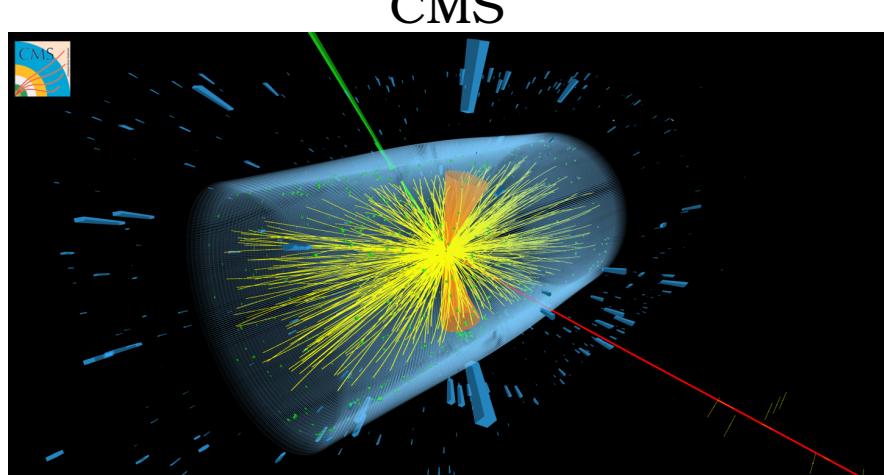
LHCb



ALICE

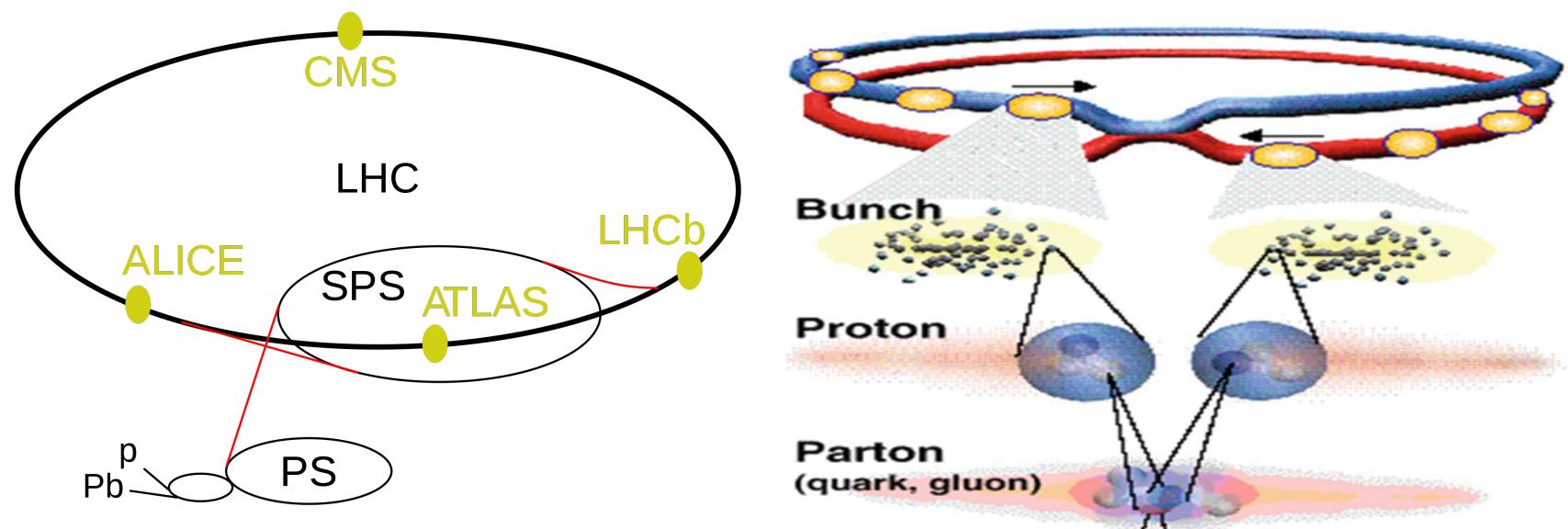


CMS



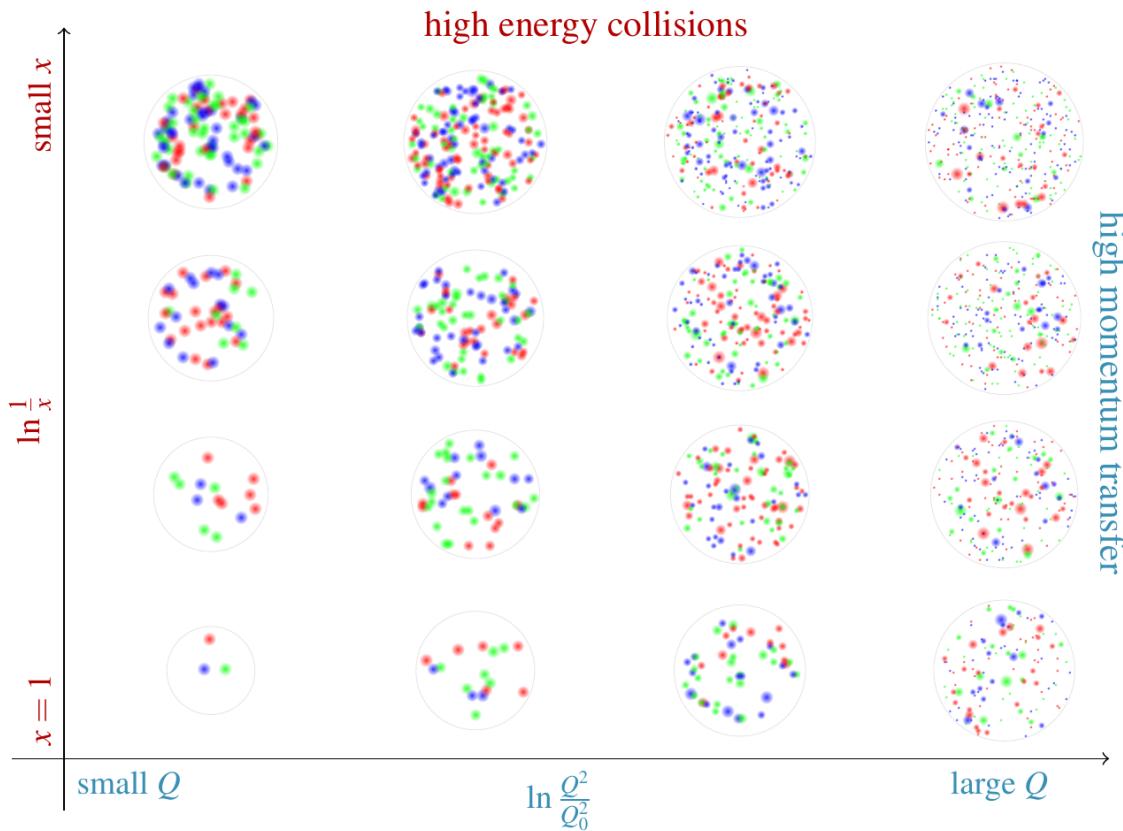
CIRCA 1 MILIARDO COLLISIONI/SEC; CIRCA 100 PETABYTE/ANNO

LHC: URTI FRA PROTONI



IL PROTONE NON È UNA PARTICELLA ELEMENTARE!

DENTRO IL PROTONE

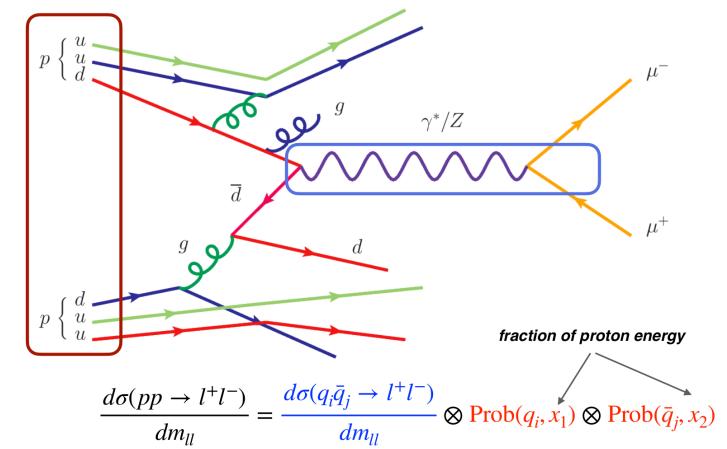


PROBABILITÀ DI PROBABILITÀ:

INFINTI² PARAMETRI

DA ESTRARRE DAI DATI

\Updownarrow



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%

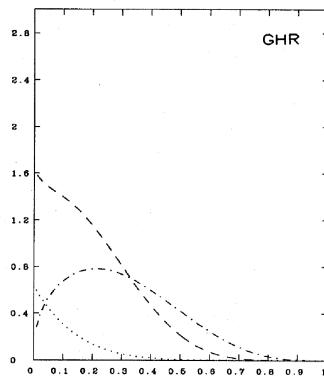


FIG. 25. Parton distributions of Glück, Hoffmann, and Reya (1982), at $Q^2=5$ GeV 2 : valence quark distribution $x[u_v(x)+d_v(x)]$ (dotted-dashed line), $xG(x)$ (dashed line), and $q_v(x)$ (dotted line).

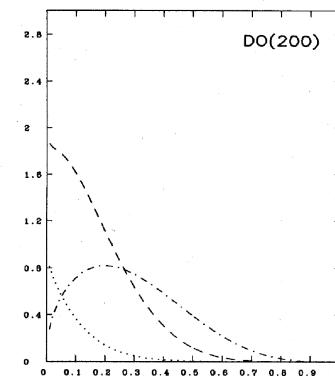


FIG. 27. "Soft-gluon" ($\Lambda=200$ MeV) parton distributions of Duke and Owens (1984) at $Q^2=5$ GeV 2 : valence quark distribution $x[u_v(x)+d_v(x)]$ (dotted-dashed line), $xG(x)$ (dashed line), and $q_v(x)$ (dotted line).

42

G. Altarelli et al. / Vector boson production

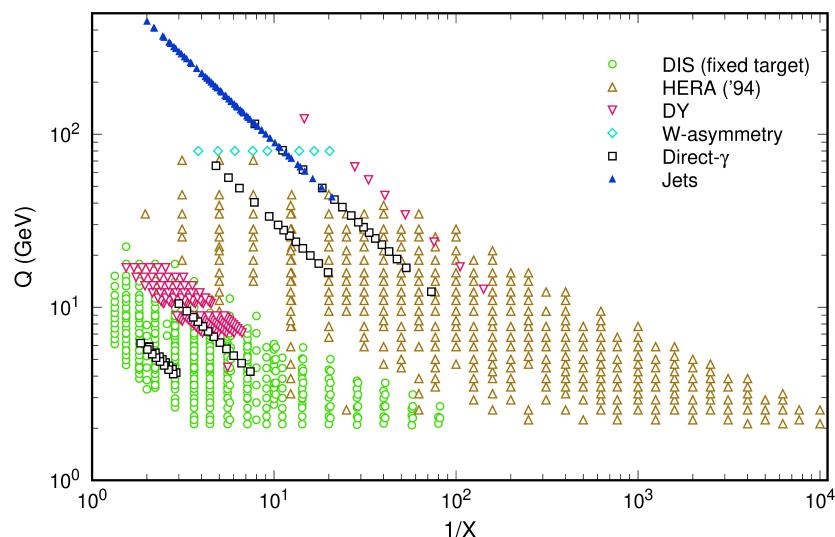
TABLE 2
Values (in nb) of the total cross sections for W^\pm and Z^0 production

\sqrt{s} (GeV)	$W^+ + W^-$ GHR	$W^+ + W^-$ DO1	$W^+ + W^-$ DO2	Z^0 GHR	Z^0 DO1	Z^0 DO2	$\frac{\sigma(W^+ + W^-)}{\sigma(Z^0)}$ GHR	$\frac{\sigma(W^+ + W^-)}{\sigma(Z^0)}$ DO1	$\frac{\sigma(W^+ + W^-)}{\sigma(Z^0)}$ 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)



PARAMETRIZZAZIONE CTEQ5 (2006)

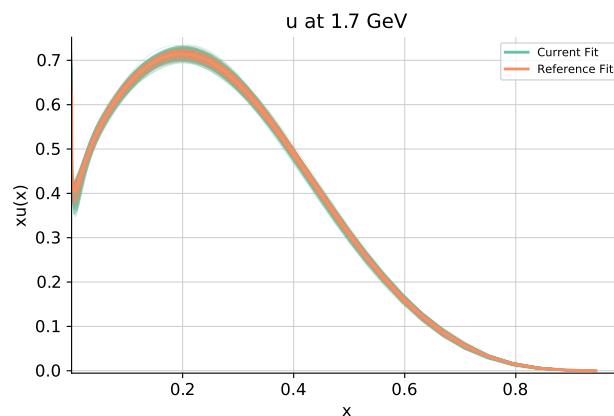
$$x f(x, Q_0) = A_0 x^{A_1} (1 - x)^{A_2} e^{A_3 x} (1 + e^{A_4} x)^{A_5}, \quad (\text{6 funzioni , 22 parametri})$$

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

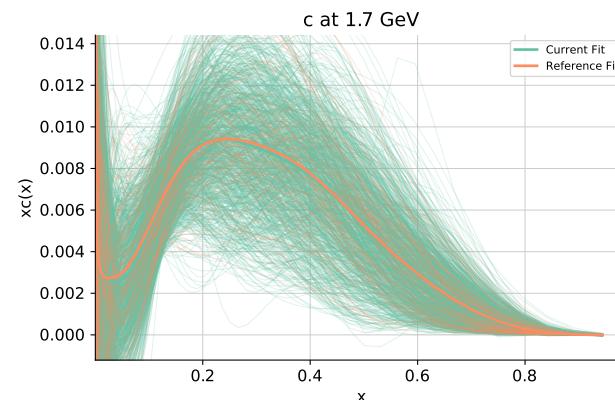
TOWARDS BIG DATA

NNPDF (2002-2017)

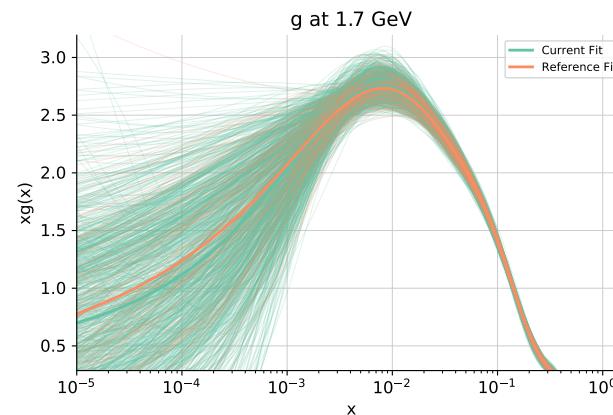
QUARK UP



QUARK CHARM



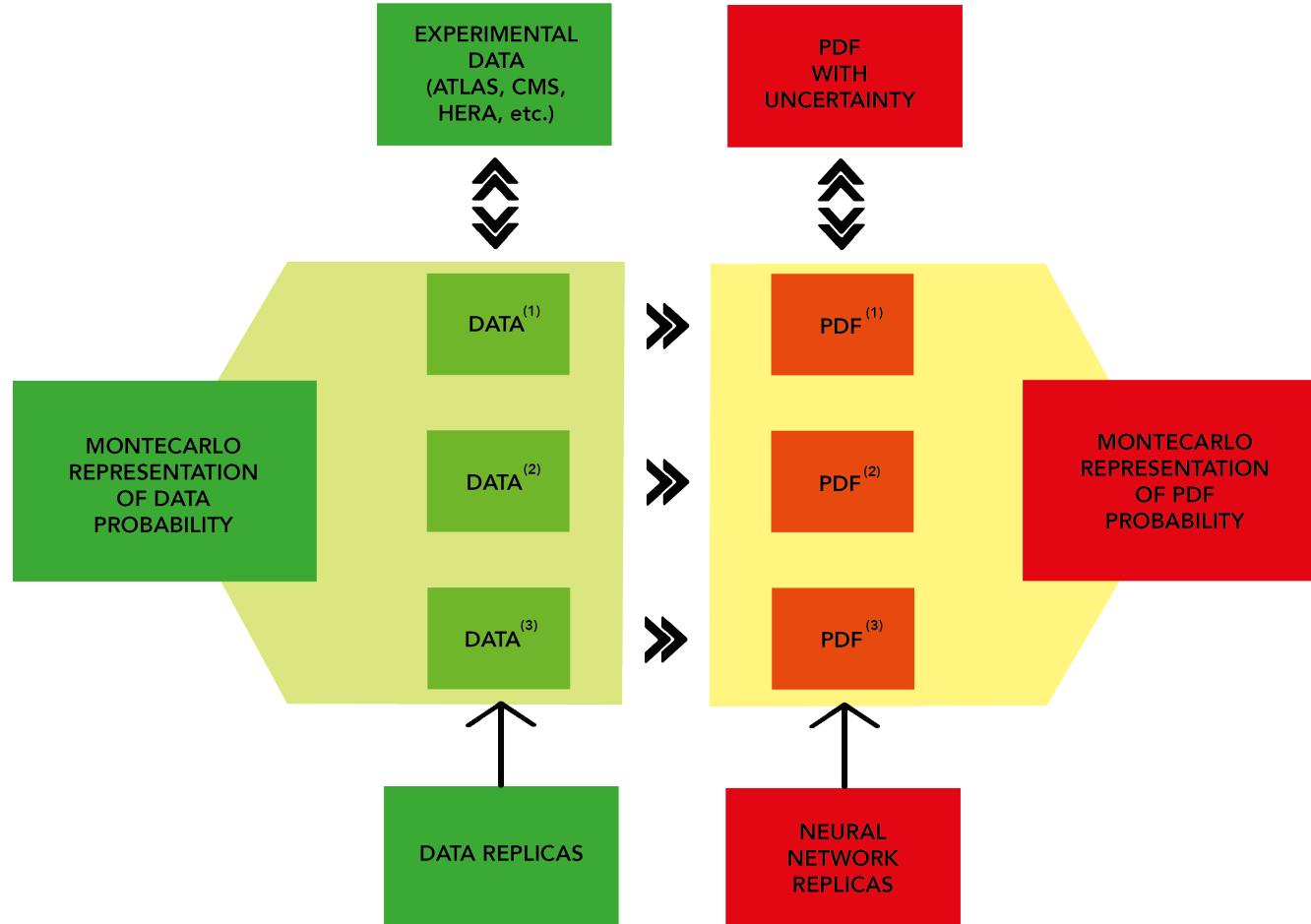
GLUONE



RAPPRESENTAZIONE MONTECARLO \Leftrightarrow DISTRIBUZIONE DI PROBABILITÀ

IL MONTECARLO FUNZIONALE

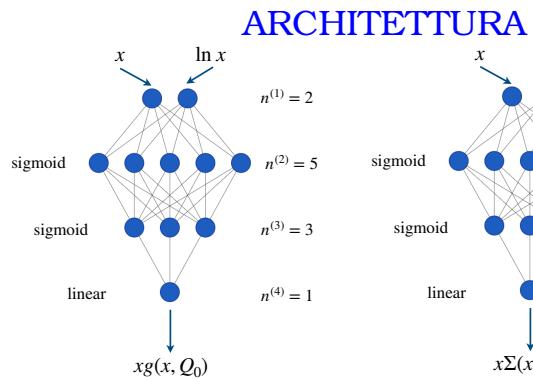
REPLICHE DI FUNZIONI \Leftrightarrow PROBABILITÀ IN UNO SPAZIO DI FUNZIONI



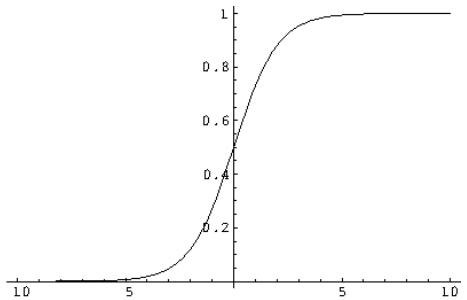
PDF OUTPUT: $f_i^{(a)}(x, \mu)$; i = up, antiup, down, antidown, strange, antistrange, charm, gluone;

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

RETI NEURALI



FUNZIONE DI ATTIVAZIONE



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

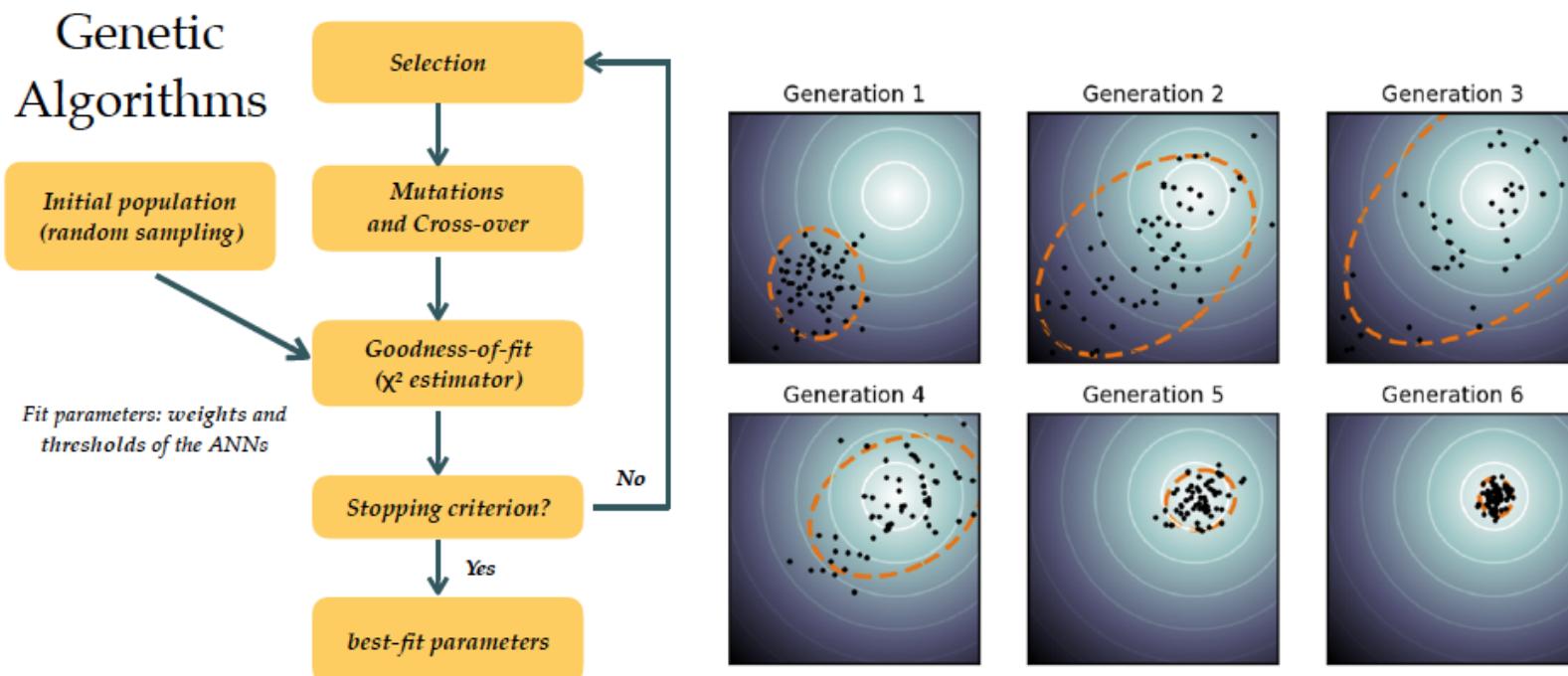
PARAMETRI

- PESI ω_{ij}
- SOGLIE θ_i
- INTERPOLANTE UNIVERSALE
- PUÒ RIPRODURRE QUALUNQUE FORMA FUNZIONALE
- COMPLESSITÀ CRESCE DURANTE L'ADDESTRAMENTO

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

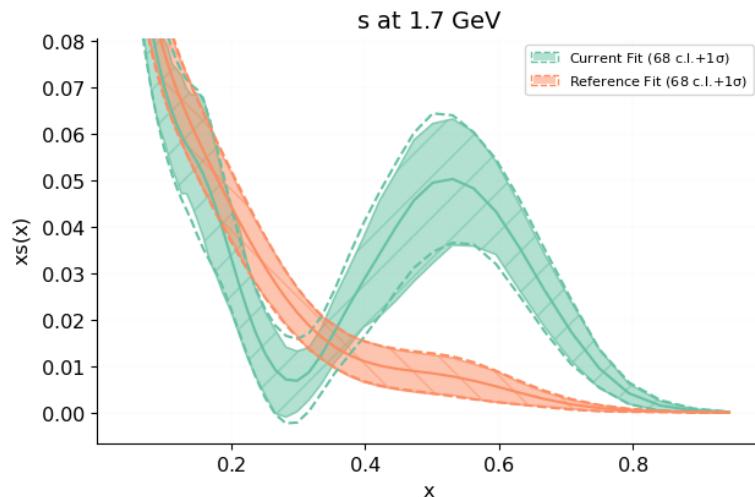
ALGORITMI GENETICI

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

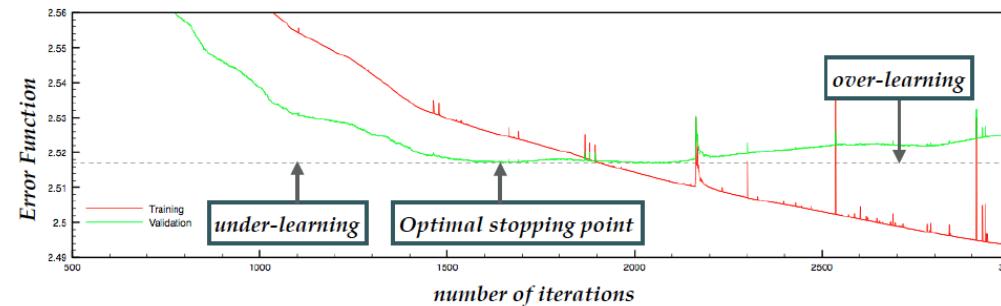


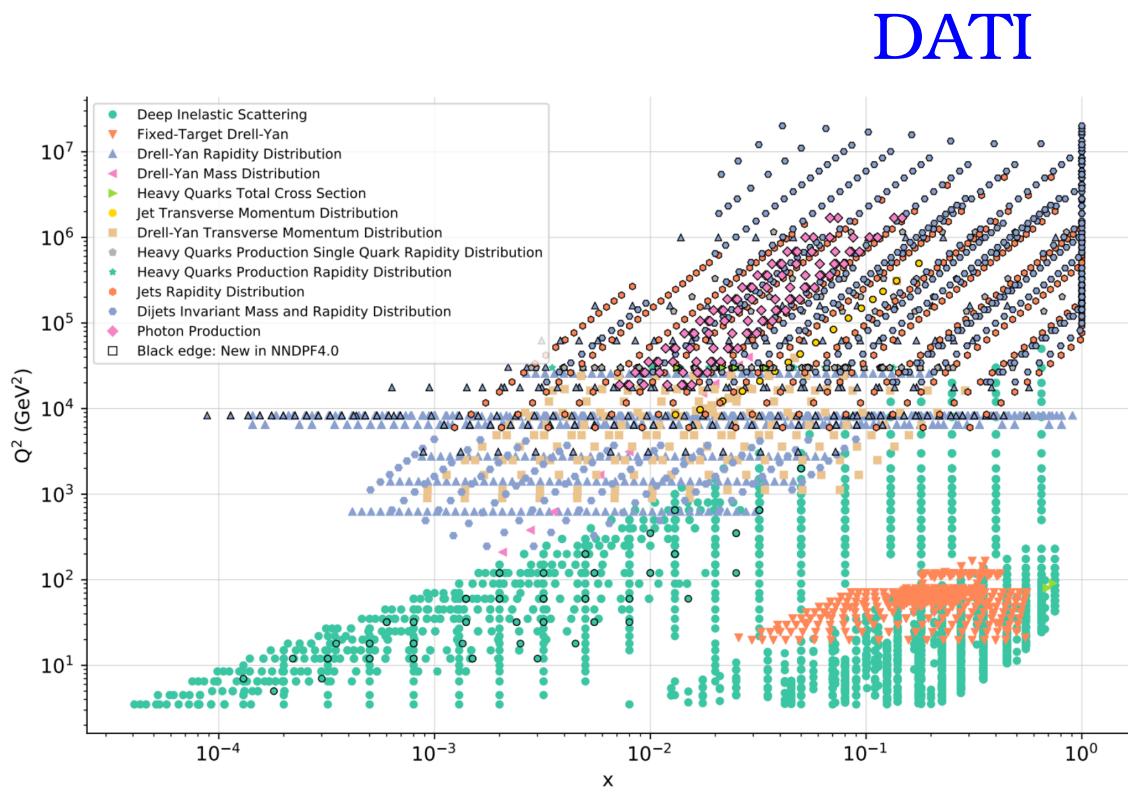
CONVALIDA INCROCIATA

- 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

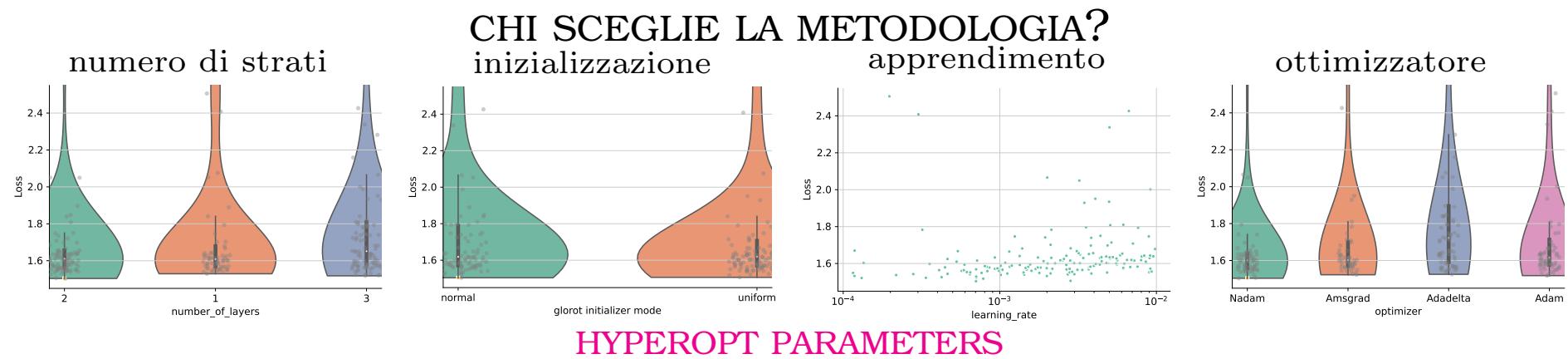




- ~ 50 ESPERIMENTI: ~ 5000 DATI SPERIMENTALI
 - 1000 REPLICHE,
 - 8 RETI NEURALI
 - ~ 20000 GENERAZIONI ALGORTIMO GENETICO
 - 80 MUTANTI PER GENERAZIONE
- $\sim 10^{13}$ PREDIZIONI TEORICHE PER UNA DETERMINAZIONE TIPICA

BIG DATA!

IPEROTTIMIZZAZIONE

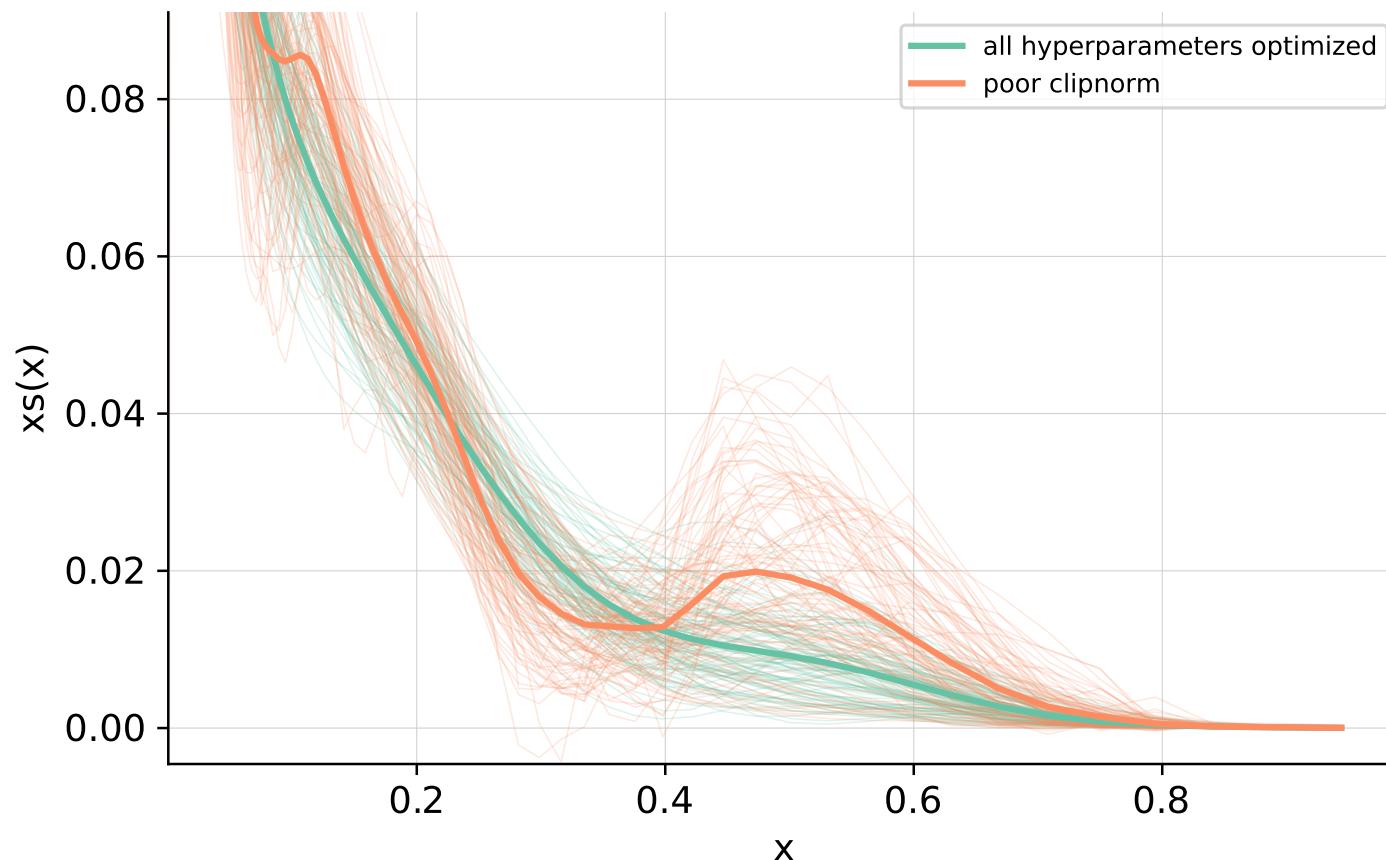


NEURAL NETWORK	FIT OPTIONS
NUMBER OF LAYERS	OPTIMIZER
SIZE OF EACH LAYER	INITIAL LEARNING RATE
DROPOUT	MAXIMUM NUMBER OF EPOCHS
ACTIVATION FUNCTIONS	STOPPING PATIENCE
INITIALIZATION FUNCTIONS	POSITIVITY& INTEGRABILITY MULTIPLIER

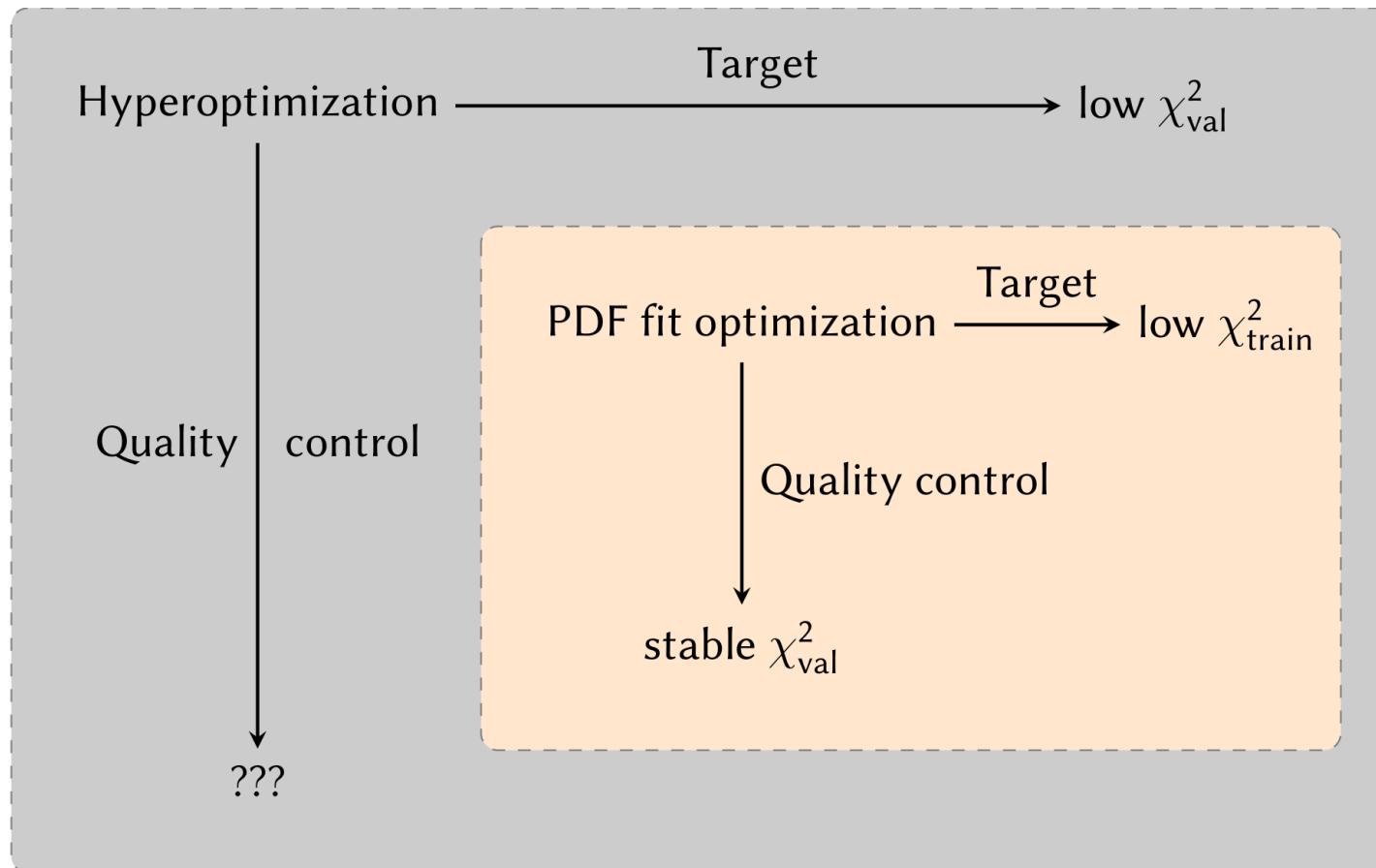
OVERFITTING

METODOLOGIA TROPPO AGGRESSIVA

s at 1.7 GeV

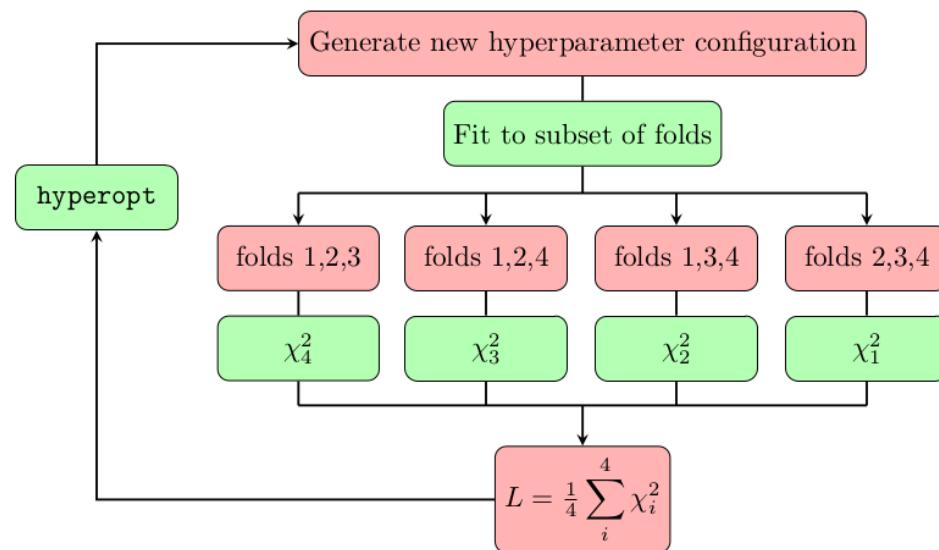


QUAL È LA METODOLOGIA OTTIMALE?



MANCA IL CRITERIO DI QUALITÀ

K-FOLDING

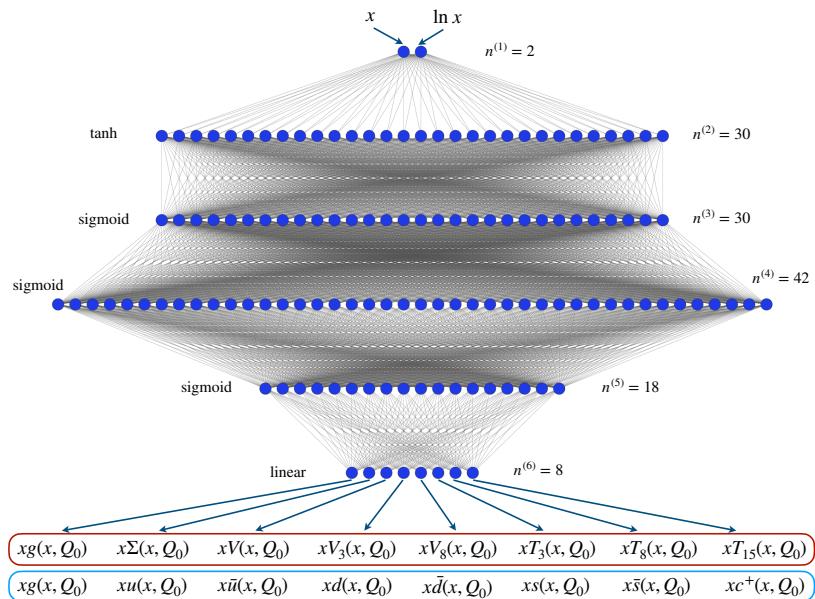


- DATI SUDDIVISI IN n GRUPPI (FOLDS)
- UN GRUPPO ESCLUSO A TURNO
- FIT MIGLIORE \Rightarrow GENERALIZZA CORRETTAMENTE AI DATI ESCLUSI

RISULTATI

15000 CICLI DI IPEROTTIMIZZAZIONE; 4 FOLDS
 FINAL NNPDF4.0 SET-UP
 AFTER HYPEROPTIMIZATION

UNA 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 2021

IL CODICE

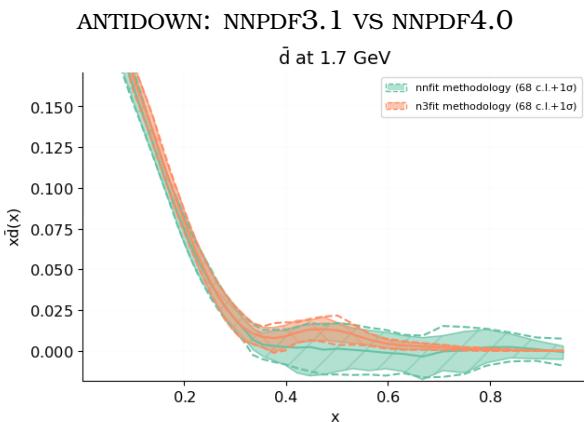
AVERAGE FITTING TIME PER REPLICA AND USE OF RESOURCES

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

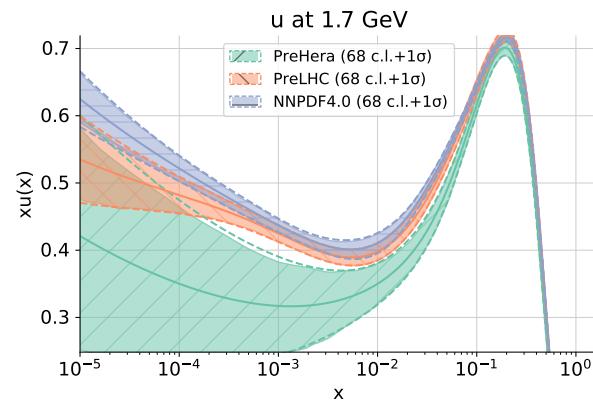
	NNPDF31 CODEBASE	NNPDF40 CODEBASE IN CPU	NNPDF40 CODEBASE IN GPU
TIME	15.2 H.	38 \pm 5 MIN.	6.6 MIN.
RAM USE	1.5 GB	6.1 GB	NA

LE PDF



IL "FUTURE TEST"

QUARK UP: DATI 1995, 2005, 2021



BIGGER DATA?

WHAT NEXT?

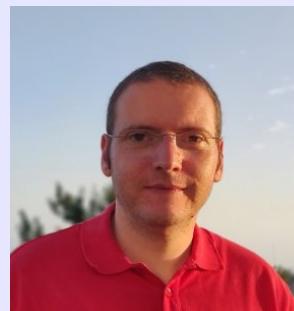
- REINFORCEMENT LEARNING
- IMPARARE LA DISTRIBUZIONE DI PROBABILITÀ
- ELEVATA PARALLELIZZAZIONE



A. Candido



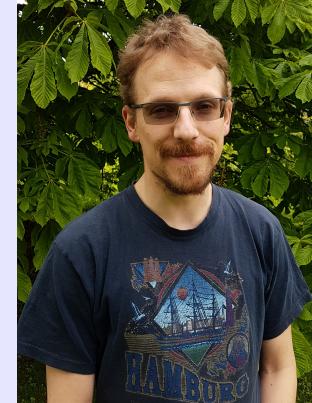
S. Carrazza



J. Cruz Martinez



F. Hekhorn



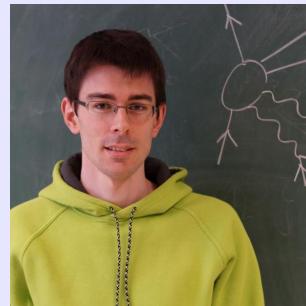
K. Kudashkin



T. Rabemananjara



C. Schwan



R. Stegeman

