## Exploiting the discovery potential of the LHC data using the Data Directed Paradigm

Shikma Bressler

Bump hunt DDP - S. Volkovitch, F. DeVito Halevi, SB [2107.11573] Symmetry DDP - M. Birman, B. Nachman, R. Sebbah, G. Sela, O. Turetz, SB [2203.07529]





# Most searches conducted following the blind analysis paradigm

- Signal region selection motivated by
  - Theoretical considerations highly motivated models are the first one to be tested
  - Final states (di-X resonance searches)
  - Topologies

# Most searches conducted following the blind analysis paradigm

- Advantage
  - Smaller chance for biassing the results
  - Best sensitivity for pre-defined signals
- Disadvantages
  - Thousands of person years are spent studying region of the data that turn out to have nothing interesting in them
  - Large portion of the data is not fully exploited



#### Example I - Resonances

	P 11 T		-	$\tau = a/a$	ь	t	č	Z/W	Н	$BSM \to SM_1 \times SM_1$			$BSM \to SM_1 \times SM_2$			$\mathrm{BSM} \to \mathrm{complex}$			
	c	μ	1	q/g	0	i	7	2/11	11	q/g	$\gamma/\pi^0{\rm `s}$	<i>b</i> ····	tZ/H	bH		$\tau qq'$	eqq'	$\mu q q'$	
e	[37, 38]	[39, 40]	[39]	ø	ø	ø	[41]	[42]	ø	ø	ø	ø	ø	ø	ø	ø	[43, 44]	ø	
μ		[37, 38]	[39]	ø	ø	ø	[41]	[42]	ø	ø	ø	ø	ø	ø	ø	ø	ø	[43, 44]	
τ			[45, 46]	ø	[47]	ø	ø	ø	ø	ø	ø	ø	ø	ø	ø	[48, 49]	ø	ø	
q/g				$\left[29, 30, 50, 51\right]$	[52]	ø	[53, 54]	[55]	ø	ø	ø	ø	ø	ø	ø	ø	ø	ø	
ь					[29, 52, 56]	[57]	[54]	[58]	[59]	ø	ø	ø	[60]	Ø	ø	ø	ø	ø	
t						[61]	ø	<b>[62]</b>	[63]	ø	ø	ø	[64]	[60]	ø	ø	ø	ø	
γ							[65, 66]	[67-69]	[68, 70]	ø	ø	ø	ø	ø	ø	ø	ø	ø	
Z/W								[71]	[71]	ø	ø	ø	ø	ø	ø	ø	ø	ø	
H									[72, 73]	[74]	ø	ø	ø	ø	ø	ø	ø	ø	
q/g										ø	ø	ø	ø	ø	ø	ø	ø	ø	
$\sim \gamma/\pi^0$ 's											[75]	ø	ø	ø	ø	ø	ø	ø	
× b												[76, 77]	ø	ø	ø	ø	ø	ø	
SN :																			
↑																			
BSN																			
:																			

#### arXiv:1907.06659

#### Mostly inclusive searches – so the data should be exploited far beyond what is seen in this table



#### Example II – LU and LFC

- LU and LFC in the SM give rise to symmetry between e's,  $\mu$ 's and  $\tau$ 's
  - Up to Yukawa interactions and phase space effect
- The discovery of an  $e/\mu$  asymmetry == New physics
  - Strong motivation to search for such asymmetries
- In practice, most of the data is not yet explored

Tested symmetry		inclusive		Obje	Object multiplicity				Decays			Topologies		
			$\tau$	j	bj	$\gamma$	Κ			W		VBF	$t\bar{t}$	
	$e/\mu$									[31-34]				
LU - U(3)	$ee/\mu\mu$ OS						[24]		[29, 30]					
	$  ee/\mu\mu$ SS													
LFC - $U(1)^{3}$	$e\mu/\mu e \text{ OS}$	$\ $ [27,28]										[27, 28]		
	$  e\mu/\mu e SS$													
	$  e^+\mu^-/e^-\mu^+$	[35]		[35]										

#### שכוז ויצמז למדע ICHEP 2022 שכוז ויצמז למדע BOLOGNA

#### Example II – LU and LFC

- LU and LFC in the SM give rise to symmetry between e's,  $\mu$ 's and  $\tau$ 's
  - Up to Yukawa interactions and phase space effect
- The discovery of an  $e/\mu$  asymmetry == New physics
- Strong motivation to search for such asymmetries
  In practice, most of the data is not yet explored

								$R_{\kappa}$	anoma	ly	
Tested symmetry		inclusive	Obj	ect mul	tiplicity	_ ]	Decays		Topo	ologies	s
Table			au j	bj	$\gamma$ K		W		VBF	$t \bar{t}$	
	$\frac{15 \text{ d VVIP}}{1 \text{ e}/\mu}$						[31-34]	1			
LU - U(3)	$ee/\mu\mu$ OS				[24]	[29,30]	[0]				
	<i>ee/µµ</i> 55										
LFC - $U(1)^3$	$e\mu/\mu e  ext{ OS} e\mu/\mu e  ext{ SS}$	[27,28]							[27, 28]		
	$\left  \begin{array}{c} e^+\mu^-/e^-\mu^+ \end{array} \right.$	[35]	[35]								



#### Searches based on the data

- Event-based anomaly detection
  - Autoencoders etc.
- Semi-supervised approaches

#### The LHC Olympics 2020

A Community Challenge for Anomaly Detection in High Energy Physics



Gregor Kasieczka (ed),<sup>1</sup> Benjamin Nachman (ed),<sup>2,3</sup> David Shih (ed),<sup>4</sup> Oz Amram,<sup>5</sup> Anders Andreassen,<sup>6</sup> Kees Benkendorfer,<sup>2,7</sup> Blaz Bortolato,<sup>8</sup> Gustaaf Brooijmans,<sup>9</sup> Florencia Canelli,<sup>10</sup> Jack H. Collins,<sup>11</sup> Biwei Dai,<sup>12</sup> Felipe F. De Freitas,<sup>13</sup> Barry M. Dillon,<sup>8,14</sup> Ioan-Mihail Dinu,<sup>5</sup> Zhongtian Dong,<sup>15</sup> Julien Donini,<sup>16</sup> Javier Duarte,<sup>17</sup> D. A. Faroughy<sup>10</sup> Julia Gonski,<sup>9</sup> Philip Harris,<sup>18</sup> Alan Kahn,<sup>9</sup> Jernej F. Kamenik,<sup>8,19</sup> Charanjit K. Khosa,<sup>20,30</sup> Patrick Komiske,<sup>21</sup> Luc Le Pottier,<sup>2,22</sup> Pablo Martín-Ramiro,<sup>2,23</sup> Andrej Matevc,<sup>8,19</sup> Eric Metodiev,<sup>21</sup> Vinicius Mikuni,<sup>10</sup> Inês Ochoa,<sup>24</sup> Sang Eon Park,<sup>18</sup> Maurizio Pierini,<sup>25</sup> Dylan Rankin,<sup>18</sup> Veronica Sanz,<sup>20,26</sup> Nilai Sarda,<sup>27</sup> Uroš Seljak,<sup>2,3,12</sup> Aleks Smolkovic,<sup>8</sup> George Stein,<sup>2,12</sup> Cristina Mantilla Suarez,<sup>5</sup> Manuel Szewc,<sup>28</sup> Jesse Thaler,<sup>21</sup> Steven Tsan,<sup>17</sup> Silviu-Marian Udrescu,<sup>18</sup> Louis Vaslin,<sup>16</sup> Jean-Roch Vlimant,<sup>29</sup> Daniel Williams,<sup>9</sup> Mikaeel Yunus<sup>18</sup>



arXiv:1807.07447

#### Searches based on the data

- Event-based anomaly detection
  - Autoencoders etc.
- Semi-supervised approaches
- Generic data/mc comparison
  - HERA  $\rightarrow$  D0 & CDF  $\rightarrow$  CMS & ATLAS
  - Sensitivity to MC mismodeling and statistics

Object Label  $p_{\rm T}$  (min) [GeV] Pseudorapidity Isolated electron 25  $|\eta| < 1.37$  or  $1.52 < |\eta| < 2.47$ P Isolated muon 25  $|\eta| < 2.7$ Isolated photon 50  $|\eta| < 1.37$  or  $1.52 < |\eta| < 2.37$ Y *b*-tagged jet 60  $|\eta| < 2.5$ 60  $|\eta| < 2.8$ Light (non-b-tagged) jet  $E_{\rm T}^{\rm miss}$ 200 Missing transverse momentum

e.g. 1µ, 3e3j, ... 10jMET, .... etc.

Data is divided to mutually exclusive classes according to object multiplicity

• A total of 704 event classes with at least 1 data event and/or 0.1 SM event prediction



8

08|s



#### Searches based on the data

- Event-based anomaly detection
  - Autoencoders etc.
- Semi-supervised approaches
- Generic data/mc comparison
- The data directed paradigm
  - Identify a property of the SM and look for regions exhibiting significant deviation from this property No MC is needed.



### The Data Directed Paradigm

- Our proposal relies solely on the data
  - Not limited by MC
- Based on two key ingredients
  - A theoretically well established property of the SM based on which deviations from the SM predictions can be searched for
  - An efficient tool that allows rapid scanning of many final states in search for such a deviation
- Complementary to ML-based method developed in an attempt to enhance the Signal/Background ratio



Tested symmetry		inclusive	$\left  \begin{array}{c} \tau \end{array} \right _{\tau}$	Object m i bi	ultip $\gamma$	licity K	 I Z	Decays W	 Topo VBF	$\frac{1}{t\bar{t}}$	s
LU - U(3)	is a WIP ee/μμ OS ee/μμ SS				/	[24]	 [29,30]	[31-34]			
LFC - U(1) <sup>3</sup>	$\begin{vmatrix} e\mu/\mu e \text{ OS} \\ e\mu/\mu e \text{ SS} \end{vmatrix}$								[27,28]		
CP	$  e^+ \mu^- / e^- \mu^+$	[35]	[	[35]							



### The Data Directed Paradigm

- Our proposal relies solely on the data
  - Not limited by MC
- Based on two key ingredients
  - A theoretically well established property of the SM based on which deviations from the SM predictions can be searched for
  - An efficient tool that allows rapid scanning of many final states in search for such a deviation
- Complementary to ML-based method developed in an attempt to enhance the Signal/Background ratio



Tested symmetry			inclusive	Object multiplicity						I	Topologies				
	is a WIP		merusive	$\tau$	j	bj	$\gamma$	Κ		Z	W		VBF	tt	
											[31-34]				
LU - U(3)	$ee/\mu\mu$ O	S						[24]		[29, 30]					
	$ $ ee/ $\mu\mu$ S	$S \parallel$													
LFC - $U(1)^{3}$	$e\mu/\mu e$ O	S	[27, 28]										[27, 28]		
	$  e\mu/\mu e S$	$S \parallel$													
CP	$ e^+\mu^-/e^- $	$\mu^+ \parallel$	[35]		[35]										



#### Example I - resonances

S. Volkovitch, F. DeVito Halevi, SB [2107.11573]

- The SM property:
  - In absence of resonances most invariant mass distributions are smoothly falling
- The tool:
  - An NN that maps invariant mass distributions to significances (q0)





#### Example I - resonances

- Looking at 1 mass bin with no signal injected
- This is the standard background-only scenario
- Using profile-likelihood tests statistics  $q_0 = Z^2$  should follow the  $\chi^2$  distributions
- This is also the case for the NN prediction

True q<sub>0</sub> Predicted q Chi2\_pdf, df = 1Chi2\_pdf, df = 1 10<sup>3</sup> 10<sup>3</sup> Entries 10<sup>5</sup> 10<sup>2</sup> 10<sup>1</sup> 10<sup>1</sup> 0 2 6 8 10 2 8 4 0 4 6 q0 true q0 predicted







#### Example I – resonances - Performance

- Significance predicted in almost no time
  Precision almost as good as that expected in blind analysis
  - The latter exploits the PLR test statistics and rely on exact knowledge of the signal and background shape
  - \*concept proved on synthetic data





Fig. 2 The difference between  $z_{\text{pred}}^{\text{max}}$  and  $z_{\text{true}}^{\text{max}}$  as a function of  $z_{\text{true}}^{\text{max}}$ . Dense regions are shown in red (roughly corresponding to the  $1\sigma$  region), while sparse regions are shown in blue.

Fig. 3 The distribution of  $z_{\text{true}}^{\max}$  (solid line) and  $z_{\text{pred}}^{\max}$  (dashed line) for samples with no signal added (blue) and for samples with a  $3\sigma$  significance signal added (orange)



#### Example II - symmetries

M. Birman, B. Nachman, R. Sebbah, G. Sela, O. Turetz, SB [2203.07529]

- The SM property:
  - Any exact or approximate symmetry of the SM:  $\underline{e/\mu}$ , CP, forward backward, ...
- The tool:
  - Symmetries allows splitting the data into two mutually exclusive datasets
    - Under the symmetry assumption they originate from the same underline distributions
  - $N_{\sigma}$  test statistics that identifies rapidly asymmetries between two datasets –

in this case the datasets are projected to histograms

$$N_{\sigma}(B,A) = \frac{1}{\sqrt{M}} \sum_{i=1}^{M} \frac{B_i - A_i}{\sqrt{A_i + B_i}}$$



## Example II - symmetries

- An "ideal" analysis
  - Background shape is perfectly known
  - Signal shape and resolution are perfectly know
  - 0 uncertainties
  - Sensitivity calculate with profile likelihood test statistics
- Analysis relaying on symmetry considerations
  - Data is split into two mutually exclusive sample one serves as a background estimate to the other
  - Background model from the data based on symmetry assumption
    - Systematic uncertainty due to available statistics in the "other" sample
  - Signal shape and resolution are perfectly know
  - Sensitivity calculate with profile likelihood test statistics







#### Example II – symmetries - Performance

- The  $N_{\sigma}$  calculation is rapid
- Can scan with no time as many matrices as one want
- Can scan sub-matrices
- Sensitivity is restored fast





#### Semi-supervised training vs the DDP

- First attempt to identify asymmetries using ML techniques
  - Demonstrate the potential of such methods but, so far, does not perform as good as the  $N_{\sigma}$  test





#### Summary

- Thousands of person hours invested so far in search for BSM physics
- Resulted in an impressive set of bounds on many BSM models
- No hints for BSM physics
- The data is far from being fully exploited New physics could easily be hidden in the already collected data
- Complementary search paradigms should be exploited
- The Data Directed Paradigm is one such possibility
  - Allows scanning rapidly many different final states and many different selection and mark those that are potentially interesting
- Concept demonstrated with two different properties of the SM
- ATLAS searches are slowly ramping up