

Where are we with  
**data citation**

Make data ~~great again~~ count

Xiaoli Chen [xiaoli.chen@cern.ch](mailto:xiaoli.chen@cern.ch)  
THOR Bootcamp Budapest  
2017. 09. 29



# Agenda

---

- Why cite data?
- How data citation works
- Principles of data citation
- Dynamic data citation
- Conclusion



---

Why cite data?



## Why cite data?

---

playing both its infrared  $\Delta v = 1$  vibrational transition (see also fig. 7) and its first allowed electronic transition  $X \rightarrow A$  in the ultraviolet. Molecular data on this electronic transition can be found in refs. [65–70]. Finally, we depict in red the sensitivity to absorption of hidden photons heavier than 11 eV onto the first three



## Why cite data?

playing both its infrared  $\Delta v = 1$  vibrational transition (see also fig. 7) and its first allowed electronic transition  $X \rightarrow A$  in the  $u$  electronic transition. Usually, we depict in r hidden photons heavy

- [65] D. M. Cooper and S. R. Langhoff, The **Journal** of Chemical Physics **74**, 1200 (1981).
- [66] C. Chackerian Jr, The **Journal** of Chemical Physics **65**, 4228 (1976).
- [67] M. Halmann and I. Laulicht, The Astrophysical **Journal** Supplement Series **12**, 307 (1966).
- [68] P. H. Krupenie, *The band spectrum of carbon monoxide*, Tech. Rep. (NATIONAL STANDARD REFERENCE DATA SYSTEM, 1966).
- [69] .
- [70] S. Tilford and J. Simmons, **Journal** of Physical and Chemical Reference Data **1**, 147 (1972).

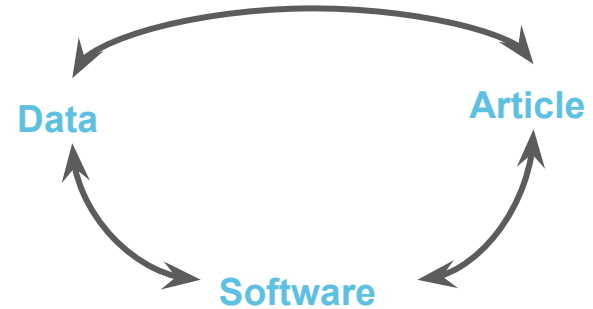
???



# Why cite data?

---

- Support proper **attribution and credit**
- Support **collaboration** and **reuse** of data
- Enable **reproducibility** of findings
- Foster **faster** and more **efficient** research **progress**, and
- Provide the **means to share** data with future researchers





---

# How data citation works

# Building a Culture of Data Citation







---

# Principles of data citation (FORCE11)



## Data citation principles (1/3)

---

1. Importance

Data is the new ~~oil~~ bacon

2. Credit and attribution

3. Evidence





# Cranmer, Kyle S.

Profile Name

Q Search

2017-09-19 14:32:11

[View Profile](#)
[Manage Profile](#)
[Manage Publications](#)
[Help](#)

## PERSONAL INFORMATION

### Personal Details (HepNames)

<b>Name</b>	Kyle S. Cranmer
<b>Current Institution</b>	New York U.
<b>E-mail</b>	<a href="mailto:cranmer@cern.ch">cranmer@cern.ch</a>
<b>Links</b>	<a href="http://theoryandpractice.org/">http://theoryandpractice.org/</a> <a href="https://www.linkedin.com/in/ky...">https://www.linkedin.com/in/ky...</a> <a href="http://twitter.com/KyleCranmer...">http://twitter.com/KyleCranmer...</a> <a href="https://github.com/cranmer">https://github.com/cranmer</a>
<b>Fields</b>	HEP-EX HEP-PH PHYSICS
<b>Experiments</b>	FNAL-E-0830 CERN-LHC-ATLAS CERN-LEP-ALEPH
<b>Identifiers</b>	BAI: <a href="#">K.S.Cranmer.1</a> INSPIRE: <a href="#">INSPIRE-00074922</a> ORCID: <a href="#">0000-0002-5769-7094</a> ARXIV: <a href="#">cranmer_k_1</a>

Period	Rank	Institution
2007	SENIOR	New York U.
2005 – 2007	PD	Brookhaven
1999 – 2005	PHD	Wisconsin U., Madison
1995 – 1999	UG	Rice U.

[Update Details](#)

## PUBLICATIONS AND OUTPUT

### Publications Datasets External

1. Data from Table 1 from: Measurement of Z-pair production in e+ e- collisions and constraints on anomalous neutral gauge couplings
2. Data from Table 1 from: Measurement of the Cross Section for open b-Quark Production in Two-Photon Interactions at LEP
3. Data from Table 1 from: Search for heavy resonances decaying to a  $W$  or  $Z$  boson and a Higgs boson in the  $q\bar{q}^{(0)}b\bar{b}$  final state in  $pp$  collisions at  $\sqrt{s} = 13$  TeV with the ATLAS detector
4. Data from Table 2 from: Search for heavy resonances decaying to a  $W$  or  $Z$  boson and a Higgs boson in the  $q\bar{q}^{(0)}b\bar{b}$  final state in  $pp$  collisions at  $\sqrt{s} = 13$  TeV with the ATLAS detector
5. Data from Table 3 from: Search for heavy resonances decaying to a  $W$  or  $Z$  boson and a Higgs boson in the  $q\bar{q}^{(0)}b\bar{b}$  final state in  $pp$  collisions at  $\sqrt{s} = 13$  TeV with the ATLAS detector
6. Data from Table 0: Final Variable Distribution from: Search for pair production of heavy vector-like quarks decaying to high- $p_{\text{transverse}}$   $W$  bosons and  $b$  quarks in the lepton-plus-jets final state in  $pp$  collisions at  $\sqrt{s}=13$  TeV with the ATLAS detector
7. Data from Table 4 from: Measurement of the Cross Section for  $W$  boson production in  $pp$  collisions at  $\sqrt{s}=13$  TeV with the ATLAS detector

### Co-Authors

[B.Mellado.1 \(13\)](#)  
[W.Quayle.1 \(11\)](#)  
[C.T.Potter.1 \(8\)](#)  
[I.Aracena.1 \(8\)](#)  
[M.Wielers.1 \(8\)](#)  
[S.L.Wu.1 \(8\)](#)  
[A.T.Watson.1 \(7\)](#)  
[B.Vachon.1 \(7\)](#)  
[C.Santamarina-Rios.1 \(7\)](#)  
[G.Louppe.1 \(7\)](#)  
[more](#)

### Subject Categories

Experiment-HEP (779)  
 Instrumentation (50)  
 Phenomenology-HEP (27)  
 Experiment-Nucl (26)

### Papers

	All papers	Single authored
<b>All papers</b>	<b>851</b>	<b>12</b>
Book	0	0
ConferencePaper	37	10
Introductory	0	0
Lectures	0	0
Published	729	4
Review	5	0
Thesis	1	1
Proceedings	0	0

### Frequent Keywords

ATLAS (687)  
 CERN LHC Coll (665)  
 experimental results (653)  
 p: scattering (521)

## STATS

### Citations Summary

851 papers found, 838 of them citeable (published or arXiv)

	Citeable papers	Published only
<b>Number of papers analyzed:</b>	838	729
<b>Number of citations:</b>	85747	82241
<b>Citations per paper (average):</b>	102.3	112.8
<b><math>h_{\text{HEP}}</math> Index [?]</b>	132	132

Breakdown of papers by citations:

	Citeable papers	Published only
Renowned papers (500+)	19	18
Famous papers (250-499)	33	31
Very well-known papers (100-249)	147	145
Well-known papers (50-99)	182	181
Known papers (10-49)	305	293
Less known papers (1-9)	118	59
Unknown papers (0)	34	2

### Citeable datasets

Number of datasets 8900

[Click here to view statistics without self-citations or RPP](#)

**Warning:** The citations count should be interpreted with great care. [Read the fine print](#)

[sign in](#) [become a supporter](#) [subscribe](#) [search](#)

jobs dating more International edition

theguardian

[UK](#) [world](#) [sport](#) [football](#) [opinion](#) [culture](#) [business](#) [lifestyle](#) [fashion](#) [environment](#) [tech](#) [travel](#) [browse all sections](#)

[home](#) > [opinion](#) > [columnists](#) > [letters](#) > [editorials](#)

**Climate change**  
Opinion

I am an Arctic researcher. Donald Trump is deleting my citations  
Victoria Herrmann

These politically motivated data deletions come at a time when the Arctic is warming twice as fast as the global average  
● Analysis: Trump signals end of US dominance in climate change battle

“As I watched more and more links turned red, I frantically combed the internet for archived versions of our country's most important polar policies.”

Tuesday 28 March 2017 11:00 BST



‘In the waning days of 2016 we were warned: save the data.’ Photograph: Andrew Stewart / SpecialStock

As an Arctic researcher, I’m used to gaps in data. Just over 1% of US Arctic waters have been surveyed to modern standards. In truth, some of the maps we use today haven’t been updated since the second world war. Navigating uncharted waters can prove difficult, but it comes with the territory of working in such a remote part of the world.

Over the past two months though, I’ve been navigating a different type of uncharted territory: the deleting of what little data we have by the Trump administration.

**Most popular**



Trump threatens to 'totally destroy' North Korea at United Nations - as it happened



Hurricane Maria: Storm grows in force to category 5 as Caribbean battered again - live



Mexico earthquake: five dead as powerful tremor rocks Mexico City



Leicester v Liverpool, Tottenham v Bamsley and more - Carabao Cup live!



Donald Trump threatens to 'totally destroy' North Korea in UN speech





## Data citation principles (2/3)

---

Proxy



Prefix



Suffix



<https://doi.org/10.5281/ZENODO.31780>



4. Unique identification

5. Access

6. Persistence

[Creative Commons Attribution-Share  
Alike 3.0 Unported](#) Alan Wilson

<https://doi.org/10.5438/55e5-t5c0>



## Data citation principles (3/3)

---

7. Specificity and Verifiability

8. Interoperability and Flexibility

Search for additional heavy neutral Higgs and gauge bosons in the ditau final state produced in  $36 \text{ fb}^{-1}$  of  $pp$  collisions at  $\sqrt{s} = 13 \text{ TeV}$  with the ATLAS detector

The ATLAS collaboration

Aaboud, Morad , Aad, Georges , Abbott, Brad ,  
Abdinov, Ovsat , Abeloos, Baptiste , Abidi, Syed  
Haider , AbouZeid, Ossama , Abraham, Nicola ,  
Abramowicz, Halina , Abreu, Henso

No Journal Information, 2017

<http://dx.doi.org/10.17182/hepdata.78402>




#### Abstract (data abstract)

CERN-LHC. A search for heavy neutral Higgs bosons and  $Z'$  bosons is performed using a data sample corresponding to an integrated luminosity of  $36.1 \text{ fb}^{-1}$  from proton-proton collisions at  $\sqrt{s} = 13 \text{ TeV}$  recorded by the ATLAS detector at the LHC during 2015 and 2016.

The  $1l1\tau_{\text{u,h}}$  channel fiducial region is defined as:

- $1e \rightarrow 1\tau_{\text{u,h}}$  or  $1\mu \rightarrow 1\tau_{\text{u,h}}$
- $p_{T\text{Lepton}} > 30 \text{ GeV}$
- $|\eta_{\text{tau}}| < 2.4$ ,  $|\eta_{\text{e}}| < 2.47$  (excluding  $1.37 < \eta_{\text{e}} < 1.57$ )



#### Table 1

Data from Figure 5A  
10.17182/hepdata.78402.v1/t2  
Observed and predicted  $m_{T\text{tot}}$  distribution in the b-veto category of the  $1l1\tau_{\text{u,h}}$  channel. Despite listing this as an exclusive final...

#### Table 2

Data from Figure 5B  
10.17182/hepdata.78402.v1/t3  
Observed and predicted  $m_{T\text{tot}}$  distribution in the b-tag category of the  $1l1\tau_{\text{u,h}}$  channel. Despite listing this as an exclusive final...

#### Table 3

Data from Figure 5C  
10.17182/hepdata.78402.v1/t4  
Observed and predicted  $m_{T\text{tot}}$  distribution in the b-veto category of the  $2\tau_{\text{u,h}}$  channel. Despite listing this as an exclusive final...

#### Table 4

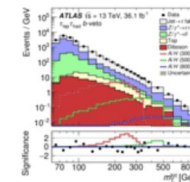
Data from Figure 5D  
10.17182/hepdata.78402.v1/t5  
Observed and predicted  $m_{T\text{tot}}$  distribution in the b-tag category of the  $2\tau_{\text{u,h}}$  channel. Despite listing this as an exclusive final...

#### Table 1 [10.17182/hepdata.78402.v1/t2](http://dx.doi.org/10.17182/hepdata.78402.v1/t2)





Observed and predicted  $m_{T\text{tot}}$  distribution in the b-veto category of the  $1l1\tau_{\text{u,h}}$  channel. Despite listing this as an exclusive final state (as there must be no b-jets), there is no explicit selection on the presence of additional light-flavour jets. Please note that the bin content is divided by the bin width in the paper figure, but not in the HepData table. In the paper, the first bin is cut off at 60 GeV for aesthetics but contains underflows down to 50 GeV as in the HepData table. The last bin includes overflows. The combined prediction for A and H bosons with masses of 300, 500 and 800 GeV and  $\tan \beta = 10$  in the hMSSM scenario are also provided.



#### cmenergies

#### observables

#### phrases

#### reactions

<b>SQRT(S)</b>	13000 GeV			
<b>LUMINOSITY</b>	$36.1 \text{ fb}^{-1}$			
<b>Channel</b>	$1l1\tau_{\text{u,h}}$			
<b>Category</b>	b-veto			
<b>Process</b>	Data	SM	$A/H(300)$	$A/H(500)$
<b><math>m_{T\text{tot}}</math> [GeV]</b>	NUM EVENTS			
50 - 70	29628	$29456.007812 \pm 356.4$	1.056363	0.023294

#### Visualize







**Principle 2: Credit and Attribution** (e.g. authors, repositories or other distributors and contributors)

**Principle 4: Unique Identifier** (e.g. DOI, Handle.). **Principle 5, 6 Access, Persistence:** A persistent identifier that provides access and metadata

Author(s), Year, Dataset Title, Data Repository or Archive, Version, Global Persistent Identifier

**Principle 7: Specificity and verification**  
(e.g. the specific version used).

Versioning or timeslice information should be supplied with any updated or dynamic dataset.



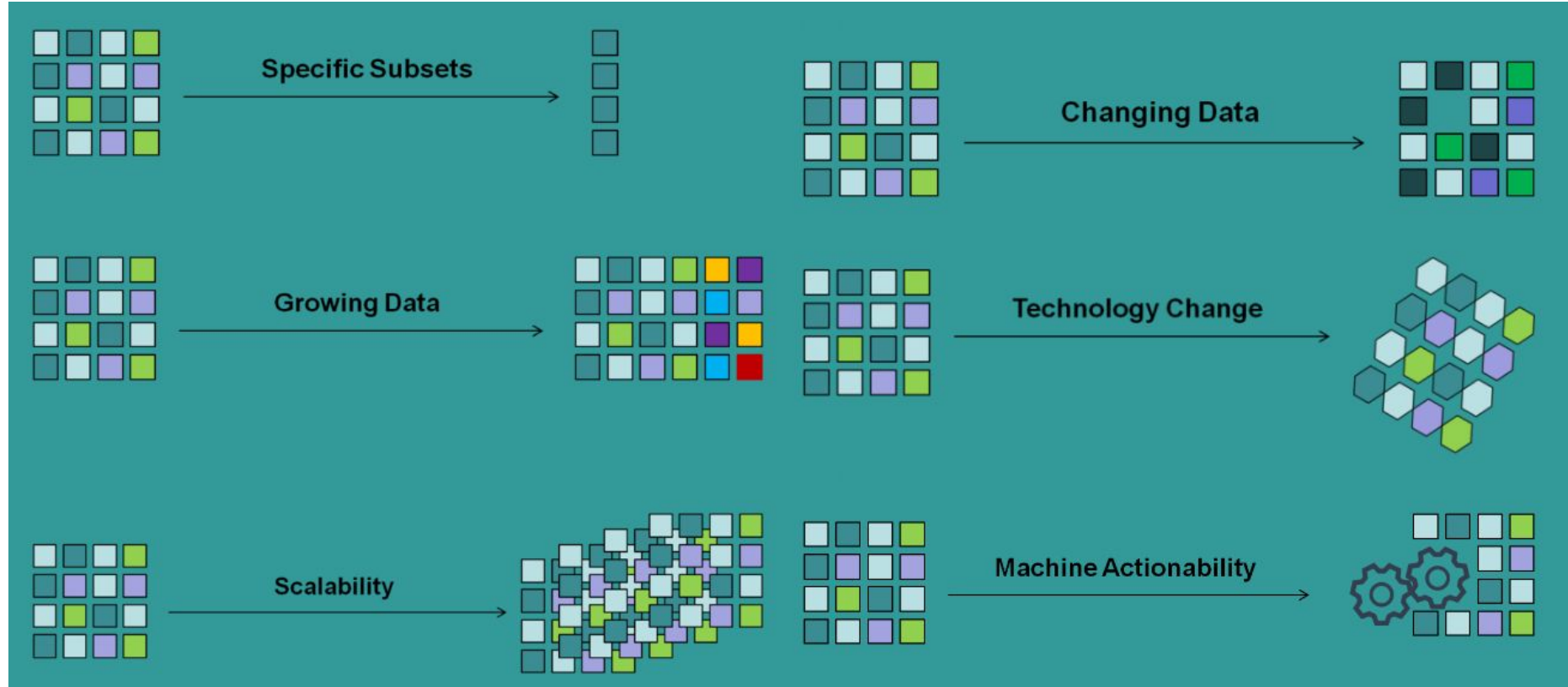


---

# Dynamic data citation



# Dynamic data citation





The WG recommends solving this challenge by:

Ensuring that data is stored in a versioned and timestamped manner.

Identifying data sets by storing and assigning persistent identifiers (PIDs) to timestamped queries that can be re-executed against the timestamped data store.



# RDA dynamic data citation recommendations

---

## DATA VERSIONING

**Keep track of meaningful versions.**

## TIMESTAMPING

**Annotate versions with timestamp.**

## QUERY STORE

**Storage for meaningful queries.**



# RDA dynamic data citation recommendations

---

## QUERY UNIQUENESS

**Only unique queries are kept on record**

## STABLE SORTING

**Defined data sorting property**

## RESULT SET VERIFICATION

**The queries e.g. data retrieval process need to be verified**



# RDA dynamic data citation recommendations

---

QUERY TIMESTAMPING

**Annotate query with timestamp**

QUERY PID

**Assign PID to query. PID resolution process = query execution process**

STORE QUERY

**Query store keeps query metadata**



# RDA dynamic data citation recommendations

---

CITATION TEXT

**Automatically generate citation text from query store metadata**

LANDING PAGE

**Human readable landing page with retrieved data and contextual information**

MACHINE ACTIONABILITY

**Provide some sort of an API**



# RDA dynamic data citation recommendations

---

TECHNOLOGY MIGRATION

**Migratable content**

MIGRATION VERIFICATION

... **self-explanatory.**





## Takeaway message

---

- Build data management requirements based on community needs
- Take advantage of the existing recommendations
- Join the on-going efforts on dynamic data citation, voice your concerns
- Let's make data first-class citizen

