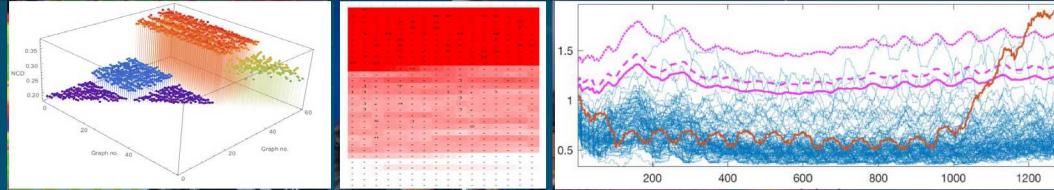




Sandia
National
Laboratories

SAND2019-2563PE

Compression Analytics for Data Science



Travis Bauer, Christina Ting, Andrew Fisher,
Rich Field, Tu-Thach Quach, Tom Brounstein,
Alex Killian, Randy Wells

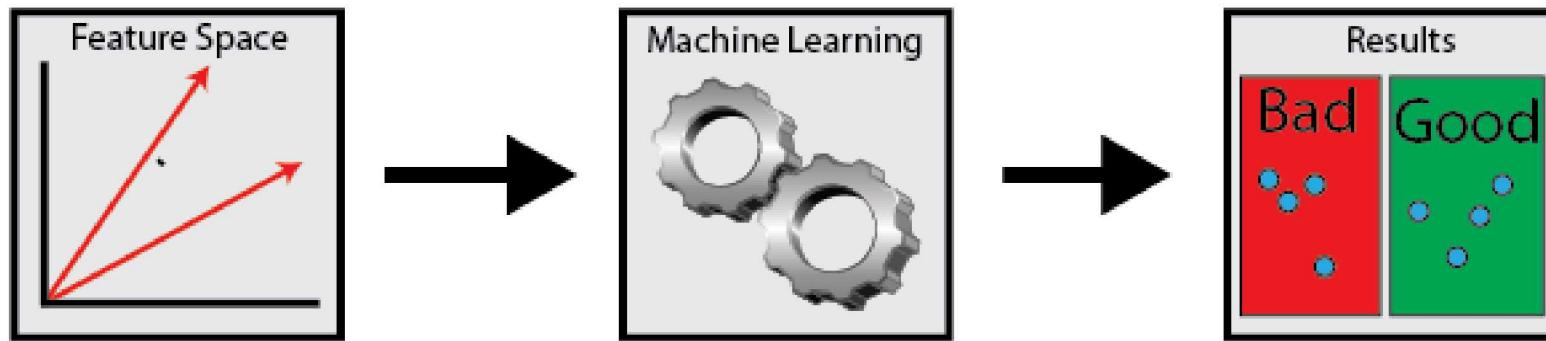


Sandia National Laboratories is a multimission laboratory managed and operated by National Technology & Engineering Solutions of Sandia, LLC, a wholly owned subsidiary of Honeywell International Inc., for the U.S. Department of Energy's National Nuclear Security Administration under contract DE-NA0003525.

Sandia National Laboratories is a multimission laboratory managed and operated by National Technology & Engineering Solutions of Sandia, LLC, a wholly owned subsidiary of Honeywell International Inc., for the U.S. Department of Energy's National Nuclear Security Administration under contract DE-NA0003525.

2

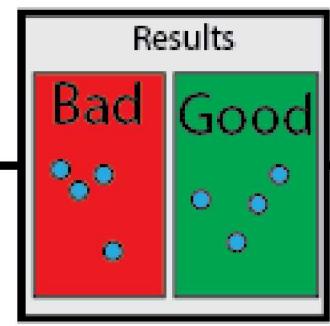
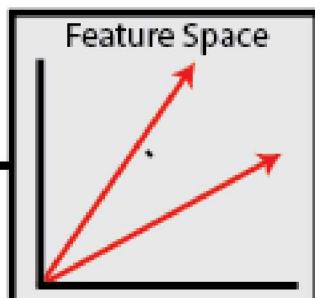
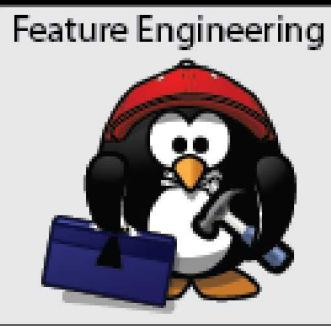
Most machine learning techniques require specific types of data.



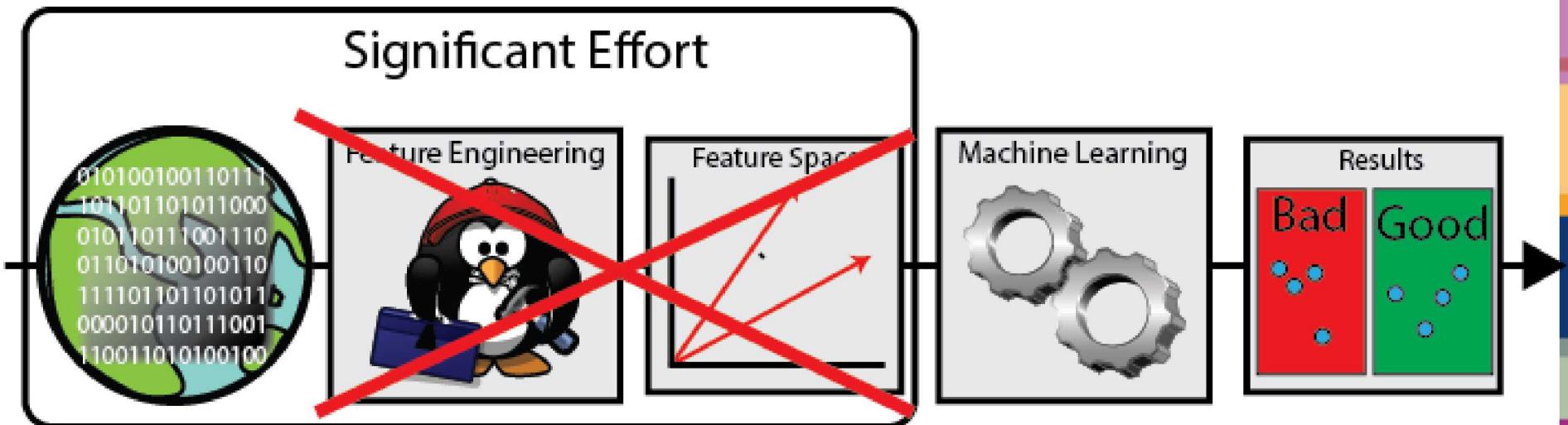
But the world doesn't come in the form of vectors and preprocessing data is a lot of work.



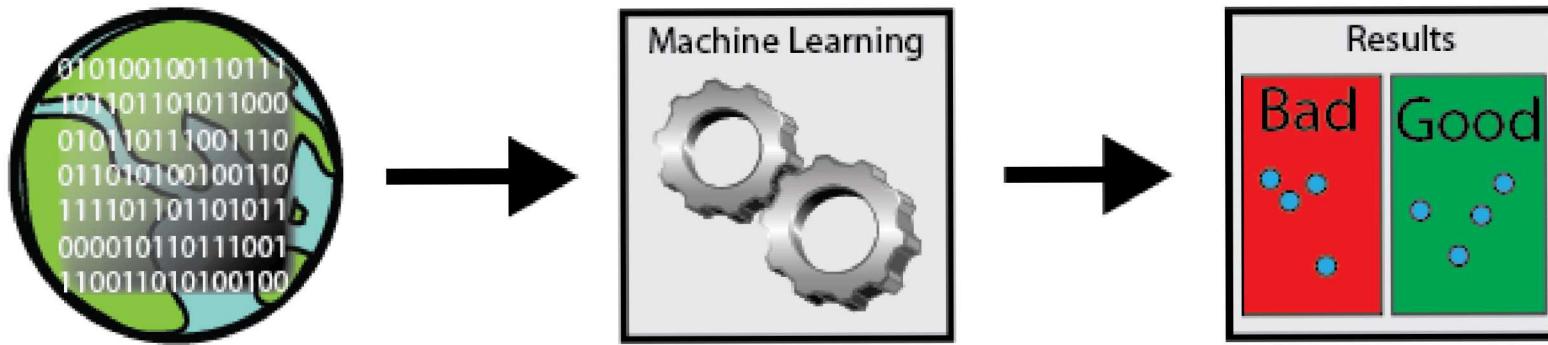
Significant Effort



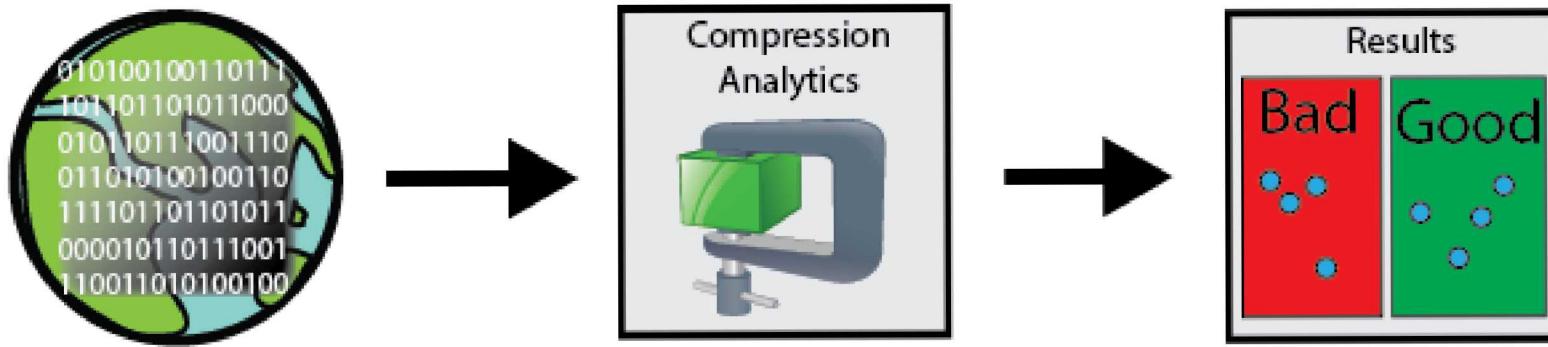
It would be nice if we could find a family of algorithms that let's us cut out the expensive feature engineering portion of our work.



It would be nice if we could find a family of algorithms that let's us cut out the expensive feature engineering portion of our work.



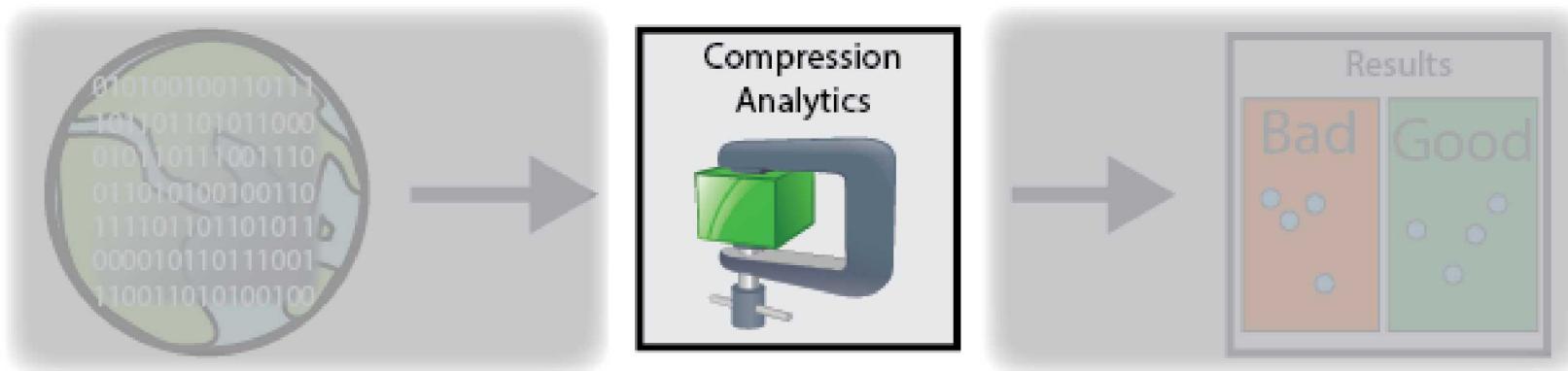
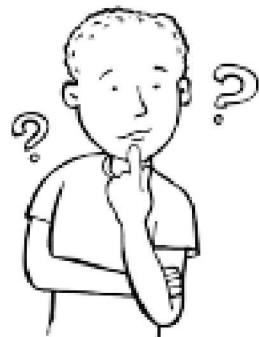
Compression-Based Analytics provides capabilities that can be used this way.



I'm going to show you how compression analytics works, where it works, and how we might make them work better



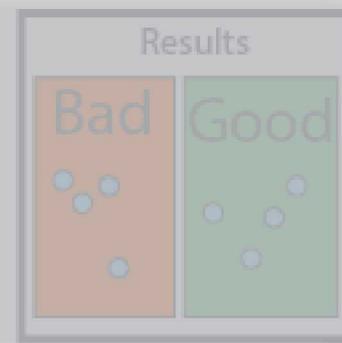
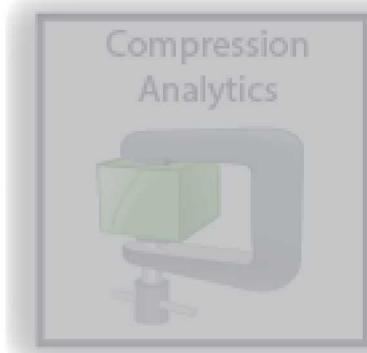
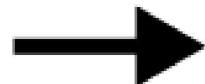
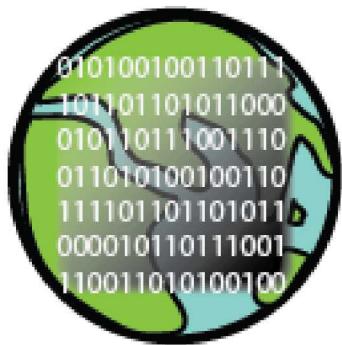
HOW compression analytics work



I'm going to show you how compression analytics works, where it works, and how we might make them work better



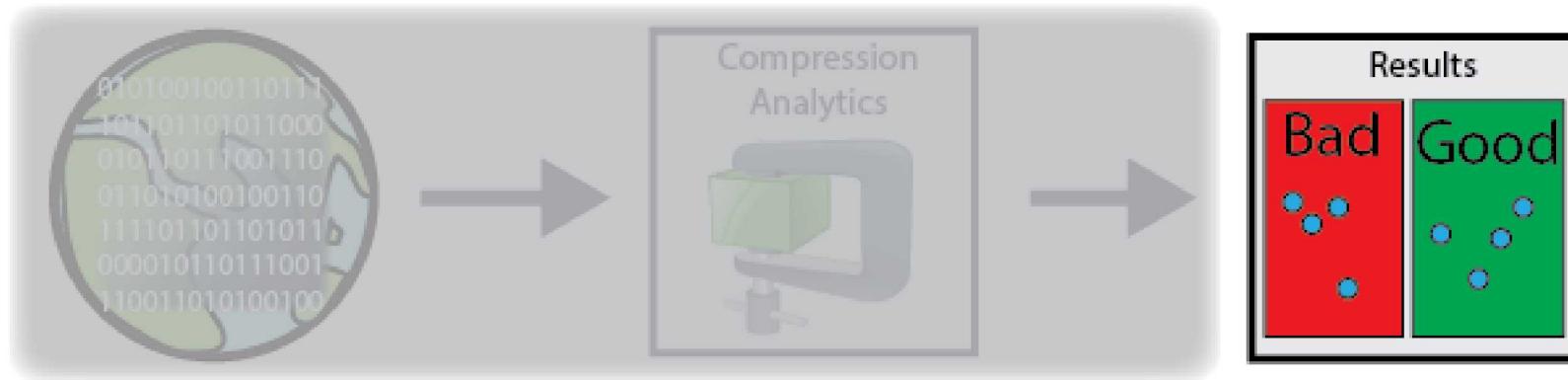
WHERE compression analytics work



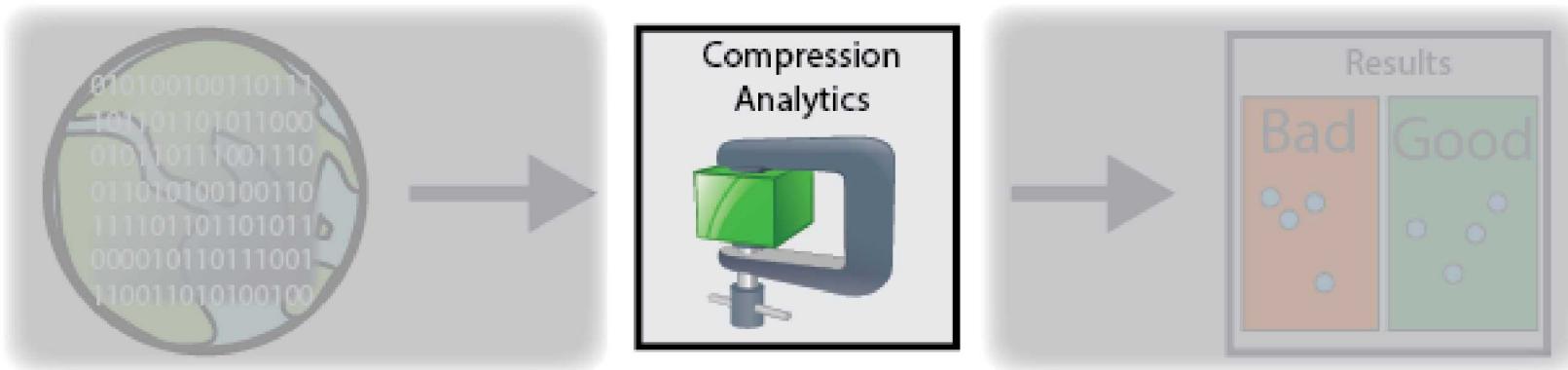
I'm going to show you how compression analytics works, where it works, and how we might make them work better



How to make compression analytics work BETTER



HOW compression analytics work



Compression analytics work by approximating the amount of shared information in data.



A String of 10K 'a's

Much Data

Almost no Information

'a'*10000

10K Truly Random Characters

Much Data

Much Information

‘TNYSTCACWIIMLZX
EWAGG...’

Compression is an upper bound on the amount of information.



A String of 10K 'a's

Much Data

Almost no Information

`Length_compressed(
'a'*10000)`

86 Bits

10K Truly Random Characters

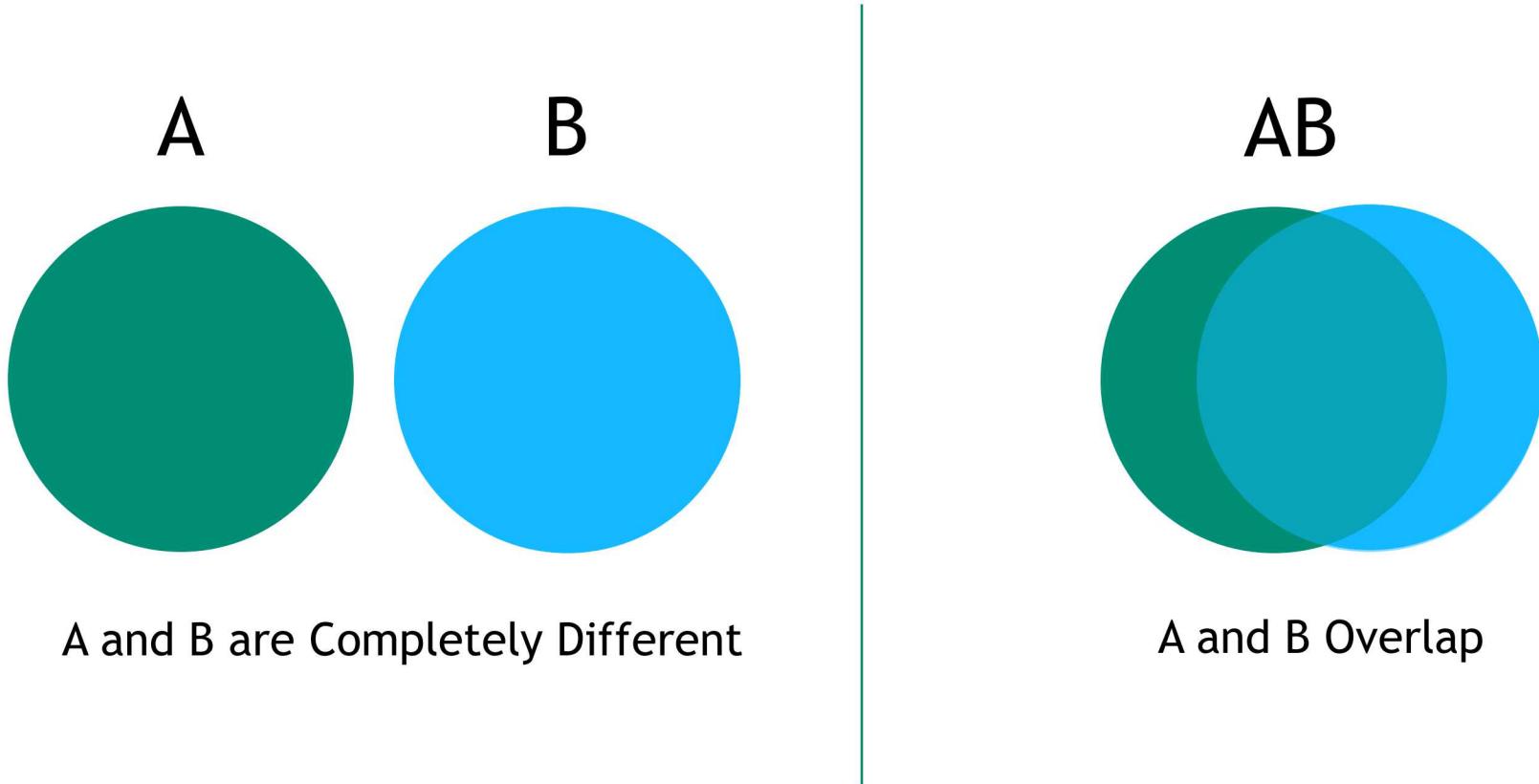
Much Data

Much Information

`Length_compressed(
'TNYSTCACWIIMLZXEWAGG...')`

68,418 Bits

We can compress two items together to approximate the amount of shared information.



$$NCD_Z(A, B) = \frac{Z(AB) - \min\{Z(A), Z(B)\}}{\max\{Z(A), Z(B)\}}$$

Li, Ming, et al. "The similarity metric." *IEEE transactions on Information Theory* 50.12 (2004): 3250-3264.

PPM/Arithmetic Coding is straight forward and a flexible implementation lends itself to analytics R&D

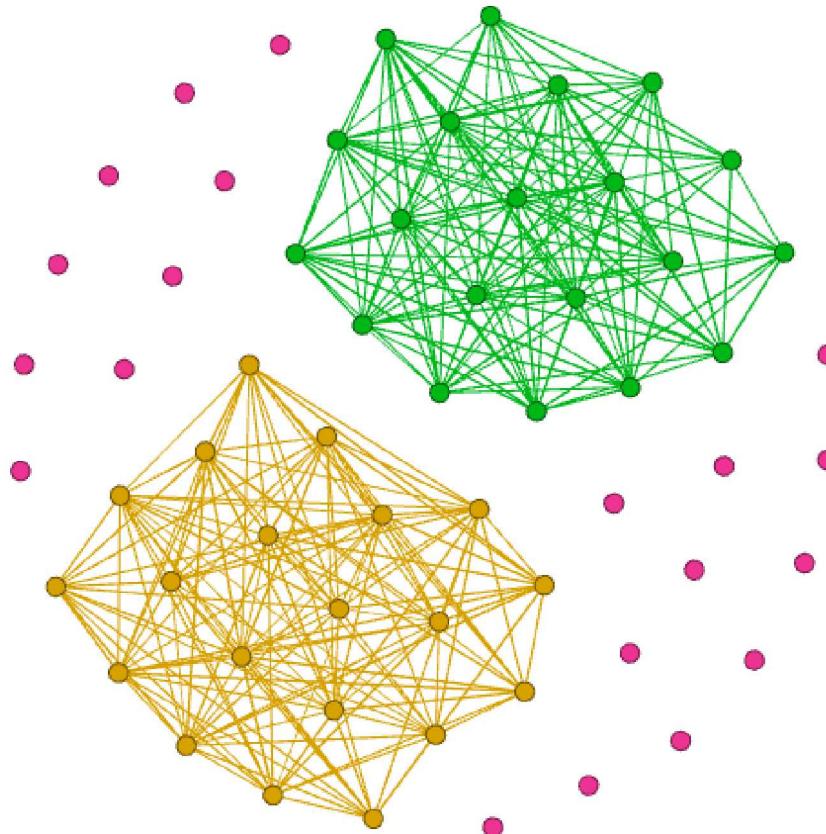


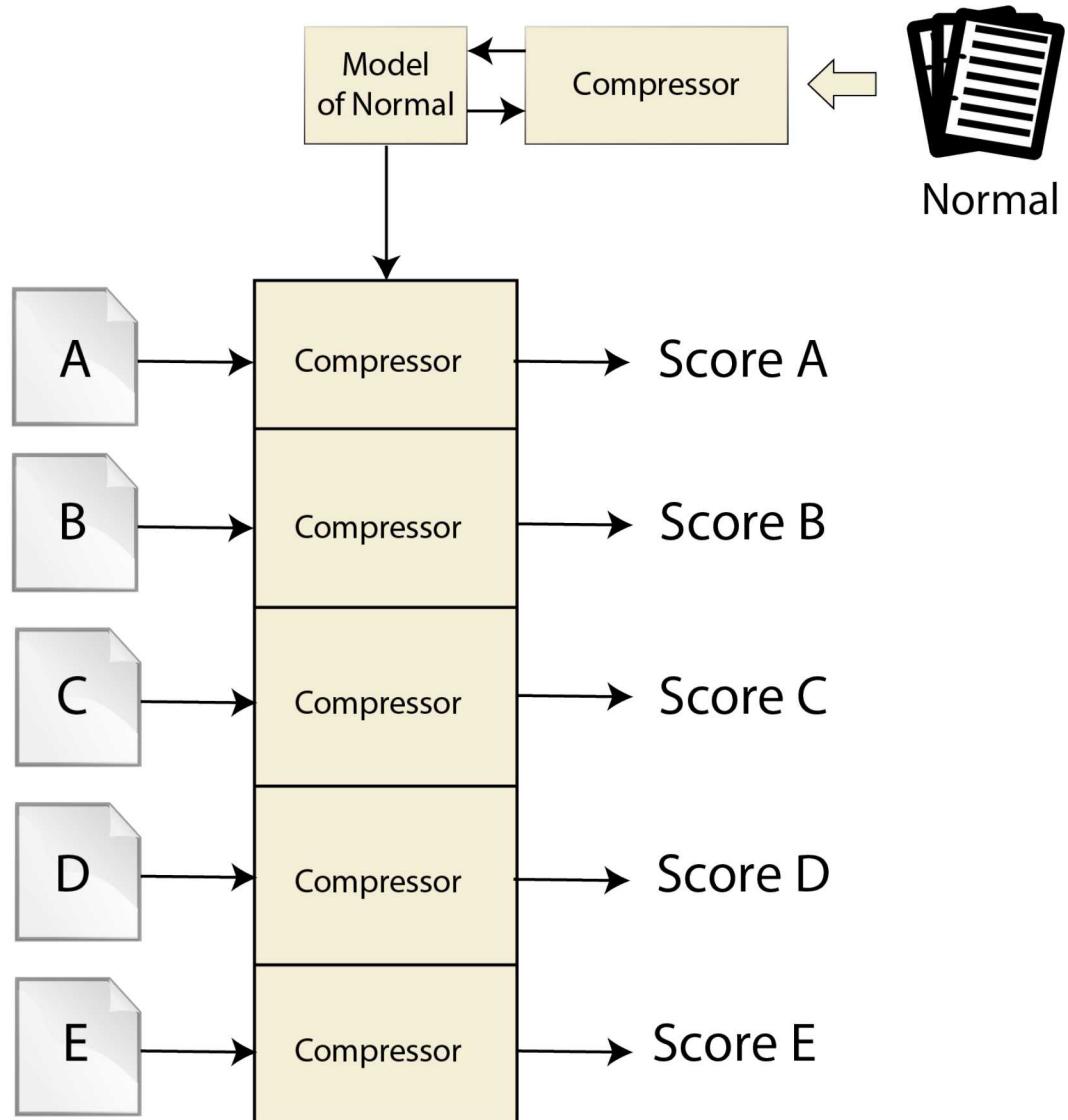
This is a question.



01010100 01101000 01101001 01110011

Distances Among Items

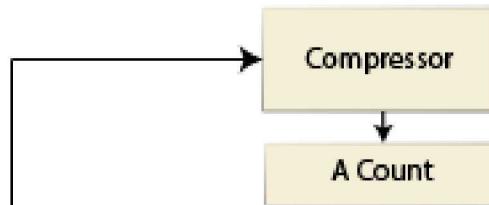




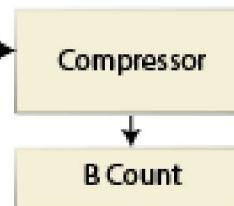


Model Application

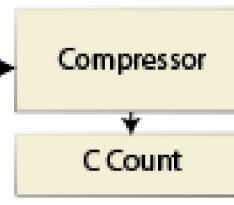
Model Building



A

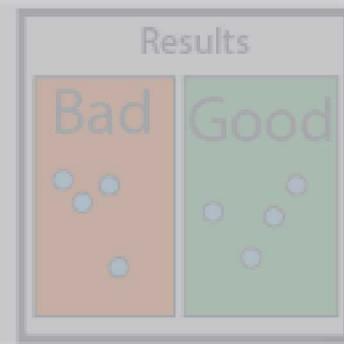
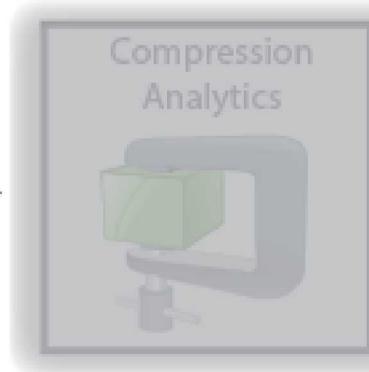
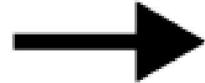
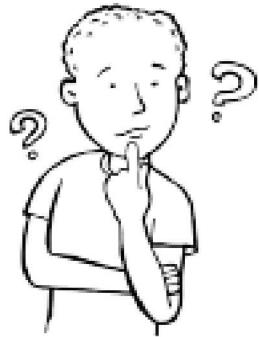


B



C

WHERE compression analytics work



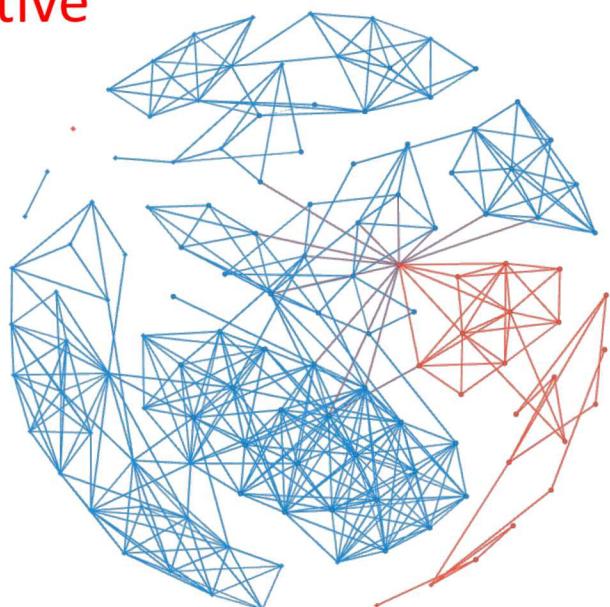
We have successfully applied compression to a variety of different data types and problems.



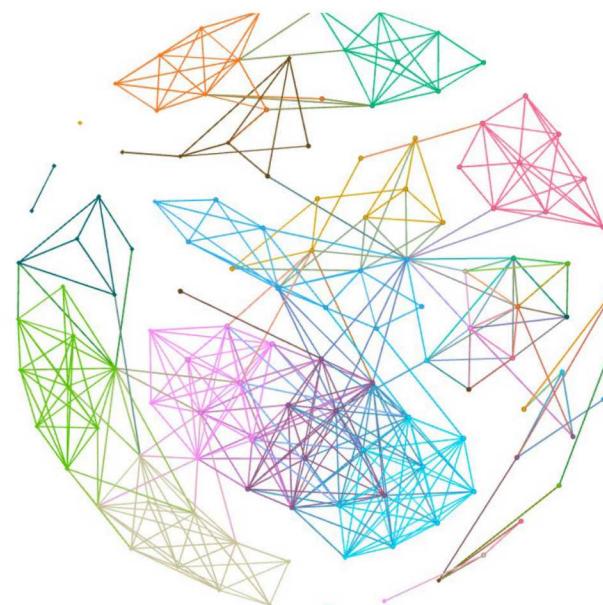
Text



Truthful vs
deceptive



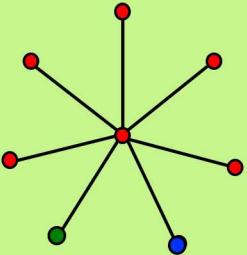
Authorship (12)



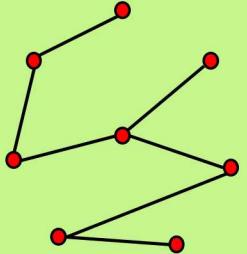
$$NCD = \frac{c(xy) - \min[c(x), c(y)]}{\max[c(x), c(y)]}$$



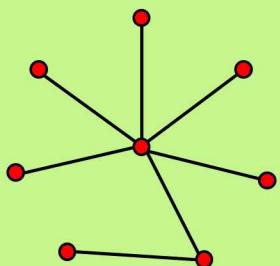
“Normal”



“Random”



“Incorrect”



```

<gexf xmlns="http://www.gexf.net/1.2draft" version="1.2">
<graph mode="static" defaultedgetype="undirected">
<attributes class="node">
  <attribute id="0" title="label" type="string" />
</attributes>

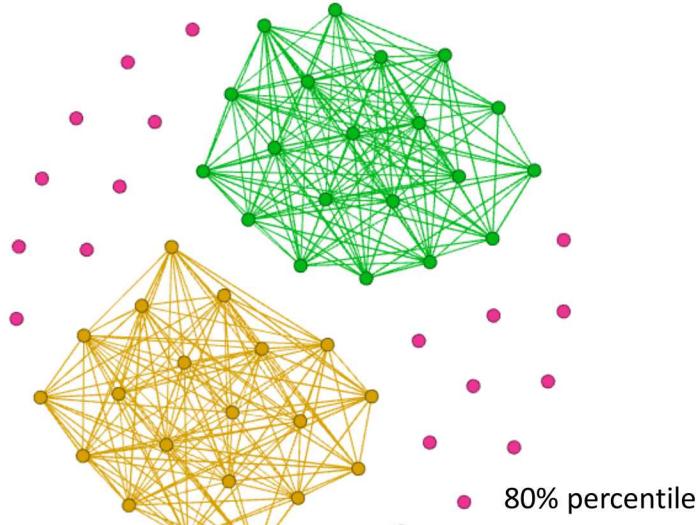
<nodes>
  <node id="1" label="root: " >
    <attvalues>
      <attvalue for="0" value="root: toolq"/>
    </attvalues>
  </node>
  <node id="2" label="name: " >
    <attvalues>
      <attvalue for="0" value="name: asqsv"/>
    </attvalues>
  </node>
  <node id="3" label="age: " >
    <attvalues>
      <attvalue for="0" value="age: ahglw"/>
    </attvalues>
  </node>
  <node id="4" label="gender: " >
    <attvalues>
      <attvalue for="0" value="gender: bjwzt"/>
    </attvalues>
  </node>
  <node id="5" label="education: " >
    <attvalues>
      <attvalue for="0" value="education: zsool"/>
    </attvalues>
  </node>
</nodes>

<edges>
  <edge id="6" source="3" target="1" type="undirected" label="1.0" weight="1.0" />
  <edge id="7" source="2" target="1" type="undirected" label="1.0" weight="1.0" />
  <edge id="8" source="4" target="1" type="undirected" label="1.0" weight="1.0" />
  <edge id="9" source="5" target="4" type="undirected" label="1.0" weight="1.0" />
</edges>
</graph>
</gexf>

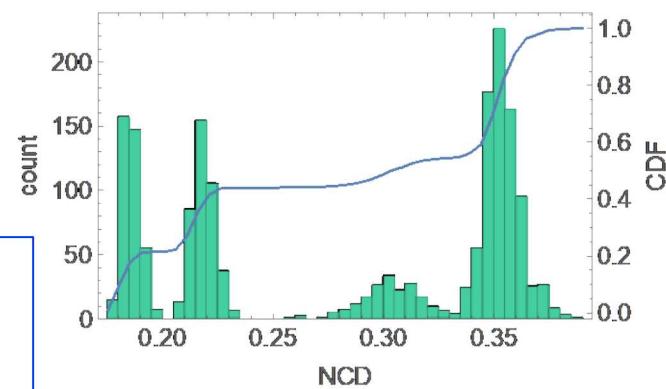
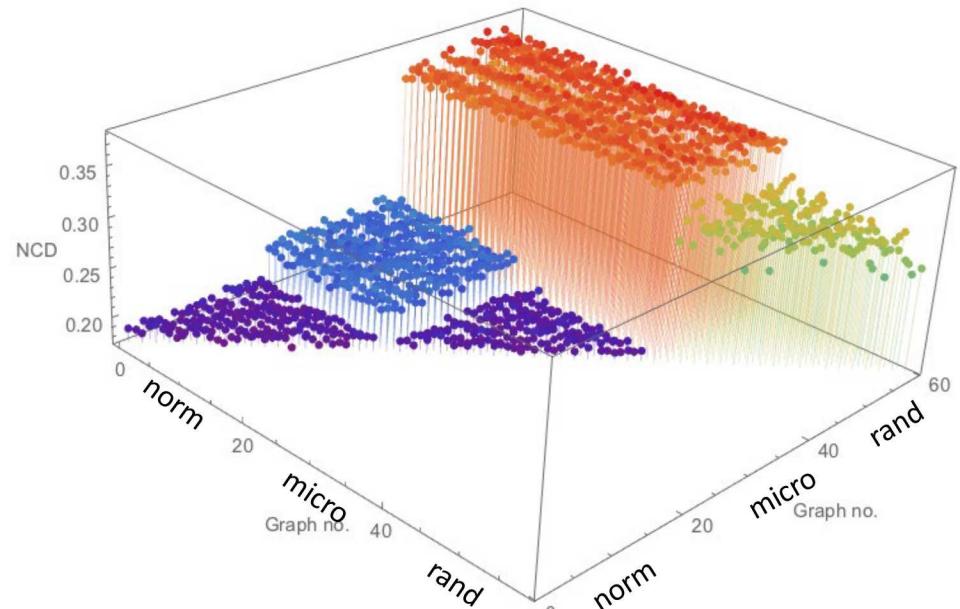
```

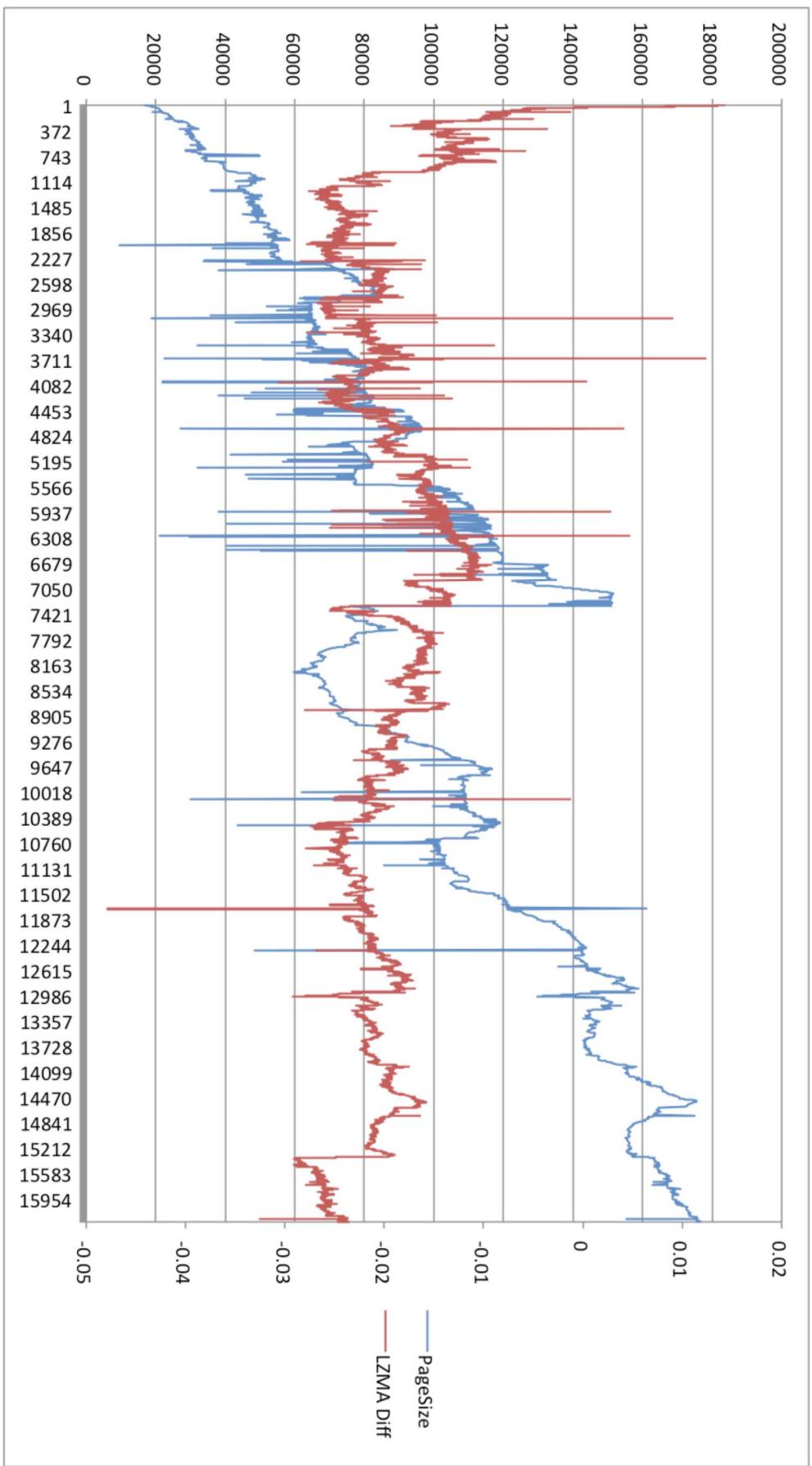


- $n = 10$
- $k = 10$
- $A = \{a, b, c, d, e\}$



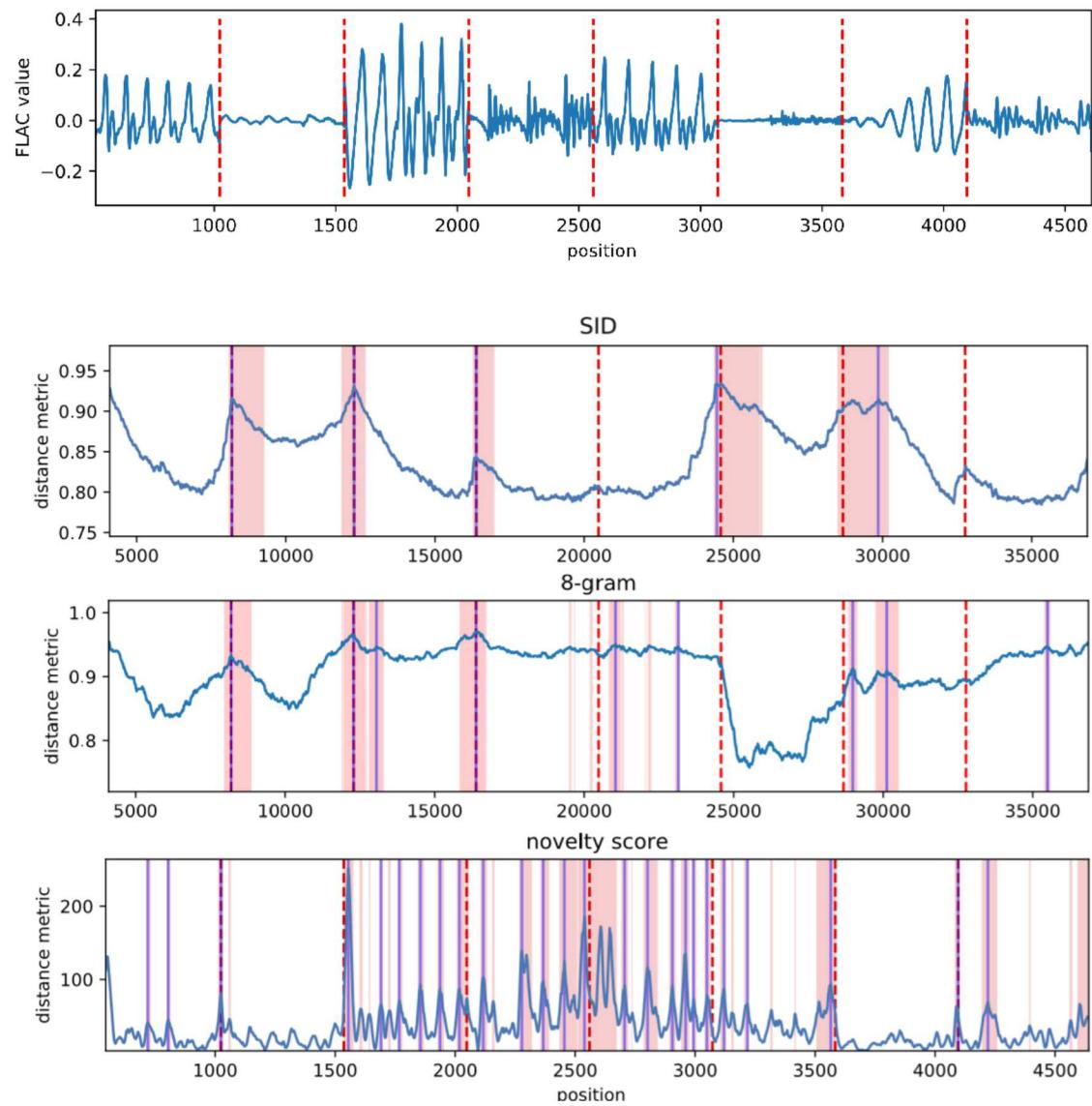
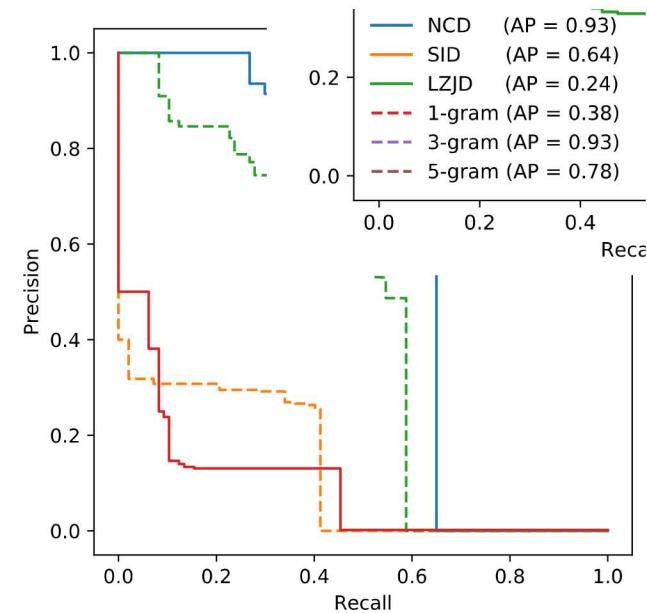
The random string, with $k = 10$ introduces variance to the NCD scores, but we can still identify a cutoff in the distribution of high similarity NCD scores

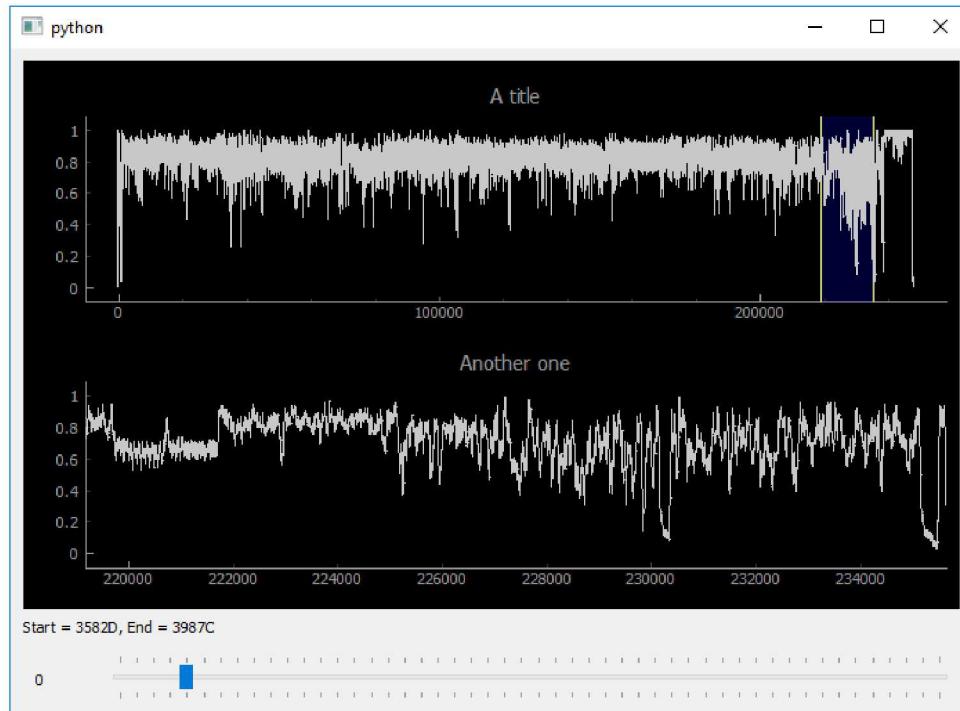
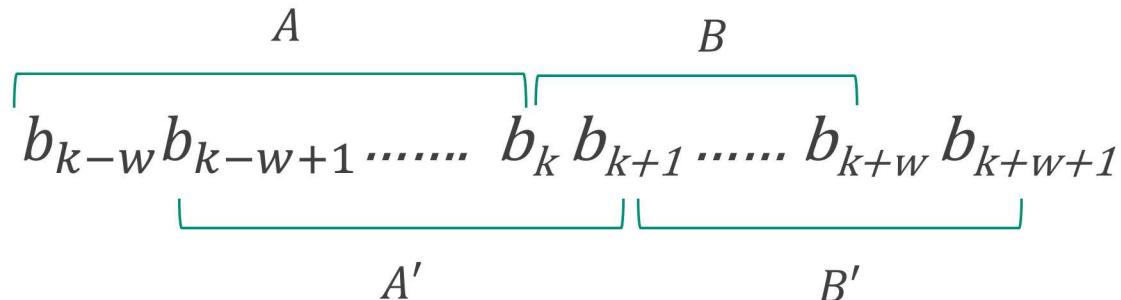






Binary

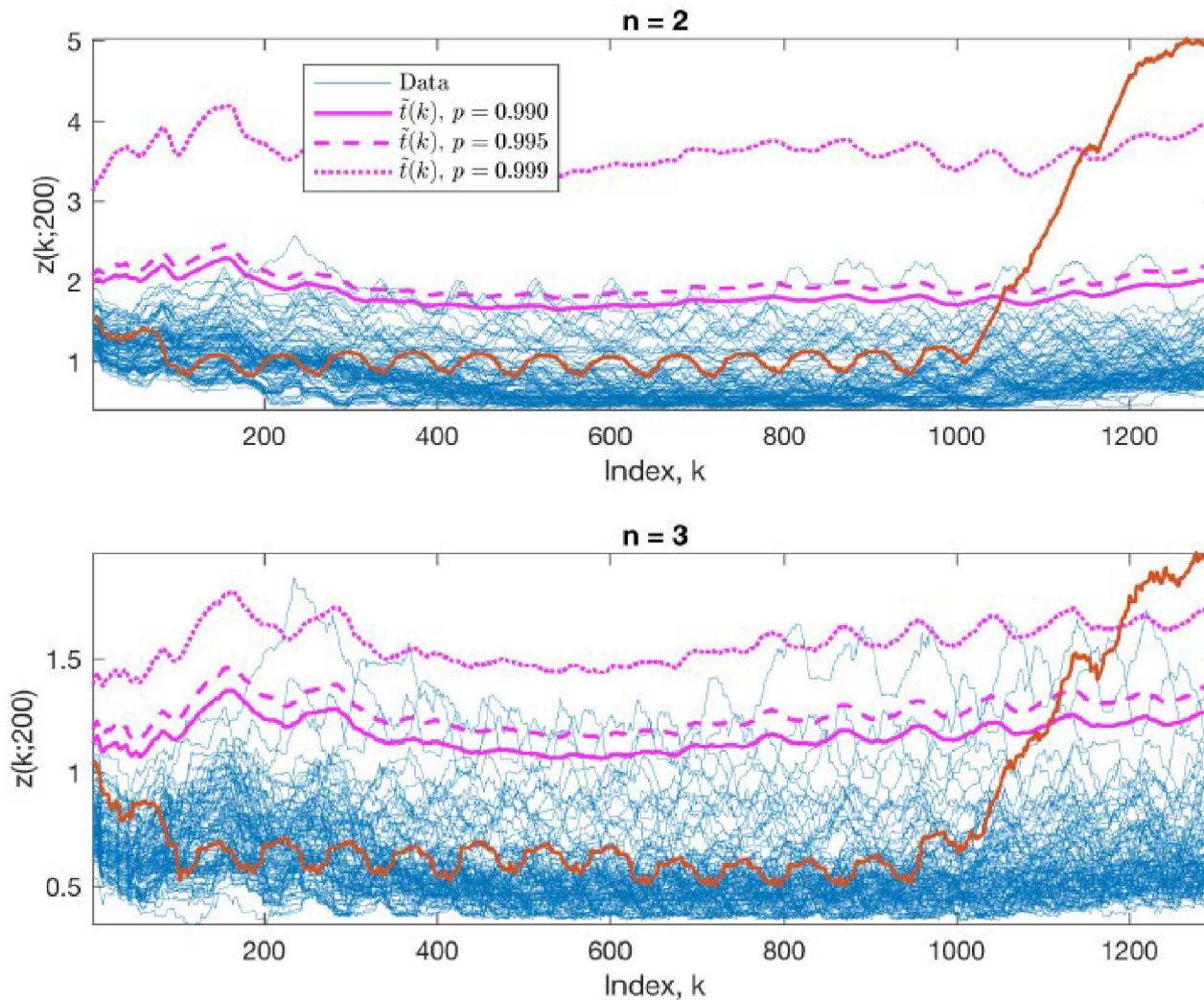




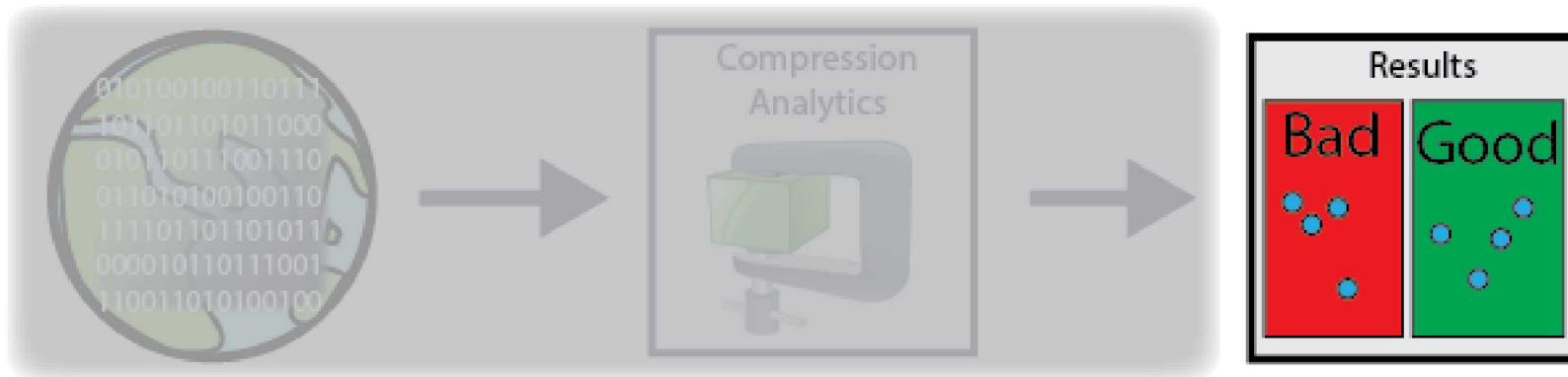
03C0	ca 22 f6 d3 9f d3 cb f9 f1 a0 9f a2 ea 0d 16 1a
03D0	33 ce 89 9d ed bf dd 4e 40 11 ec 30 24 91 ed e1 N 8 .. 0 \$..
03E0	d8 1a 93 d4 72 54 a9 65 7f e0 1f 98 c3 1b 55 0a r T .. e .. . 0 ..
03F0	96 89 ec 1f d1 66 c0 64 d6 0f a0 b7 a3 bf 66 07 f .. d .. . f ..
0400	82 25 b2 4d ce 2a 6a 15 f0 69 fc a2 57 f9 b1 11 M .. * j .. d .. W ..
0410	c7 b1 2c 5b 41 1e b3 05 8e 69 bf d7 69 0a d1 0e A i ..
0420	03 05 32 aa 6b 02 a6 9f 78 46 e2 1e f8 83 40 45 2 .. k .. . x F .. . 0 E ..
0430	c0 aa f7 ae 5c 01 3b dd 46 c9 05 65 67 04 of bf \ .. r .. F .. e g ..
0440	e6 33 f7 11 72 c9 27 fb 17 8c 88 4d e0 fd e5 c6 3 .. . r .. ' .. . M ..
0450	44 of 51 2d c4 8b 5e fe 10 90 fe ff 00 00 00 00 Q .. . ^ .. .
0460	10 90 2e 00 fo 6f 01 00 00 00 00 0a 06 00 00 o
0470	00 80 fe ff 00 00 00 00 80 2e 00 10 10 00 00
0480	00 00 00 00 0a 04 00 00 ac 6a fc ff 00 00 00 00 j
0490	ac 6a 2c 00 54 15 02 00 00 00 00 0a 06 00 00 j .. . T
04A0	7b 6a fc ff 00 00 00 00 7b 6a 2c 00 31 00 00 00 j .. . { j .. , 1 ..
04B0	00 00 00 00 0a 04 00 00 1a 6a fc ff 00 00 00 00 j
04C0	1a 6a 2c 00 61 00 00 00 00 00 00 0a 06 00 00 j .. , a
04D0	19 6a fc ff 00 00 00 00 19 6a 2c 00 01 00 00 00 j j
04E0	00 00 00 00 0a 04 00 00 00 00 f6 ff 00 00 00 00
04F0	00 00 26 19 6a 06 00 00 00 00 0a 06 00 00 00 & .. j
0500	00 00 f4 ff 00 00 00 00 24 00 00 00 02 00 \$
0510	00 00 00 00 0a 03 00 00 00 f2 ff 00 00 00 00
0520	00 00 22 00 00 00 02 00 00 00 00 0a 01 00 00 "
0530	00 00 d5 ff 00 00 00 00 05 00 00 00 1d 00
0540	00 00 00 00 0a 03 00 00 00 d2 ff 00 00 00 00
0550	00 00 02 00 00 00 02 00 01 00 00 00 00 00 00 00
0560	00 00 d0 ff 00 00 00 00 00 00 00 00 00 02 00
0570	00 00 00 00 00 02 00 00 00 d4 ff 00 00 00 00
0580	00 00 04 00 00 00 01 00 01 00 00 00 00 00 00 00
0590	00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00
05A0	00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00



Network Traffic



How to make compression analytics work BETTER

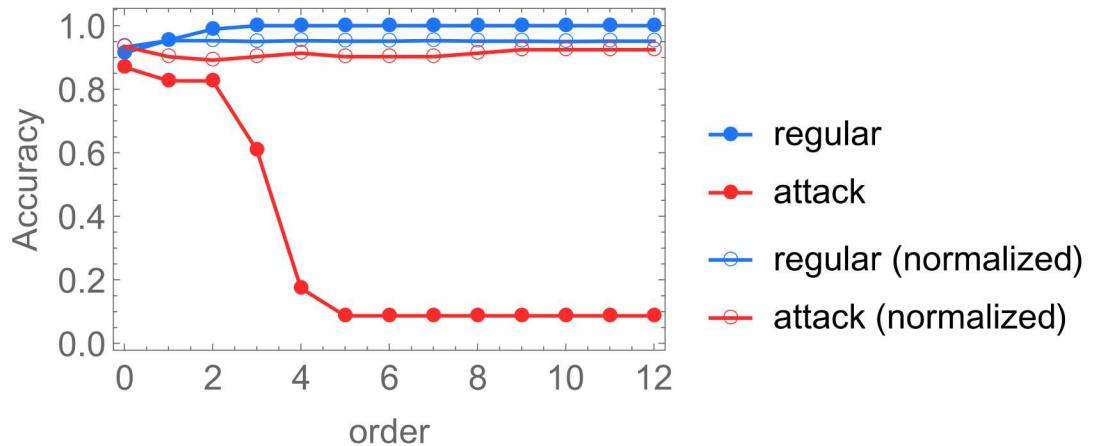
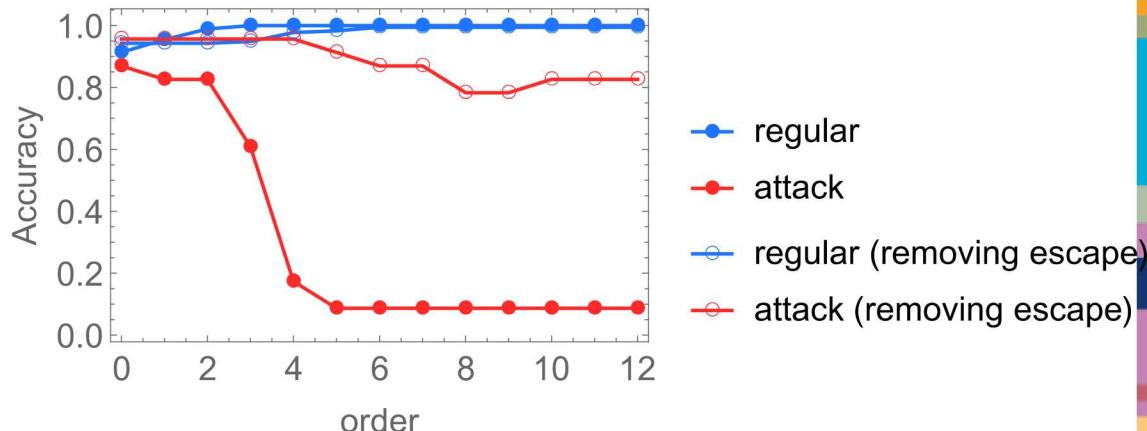
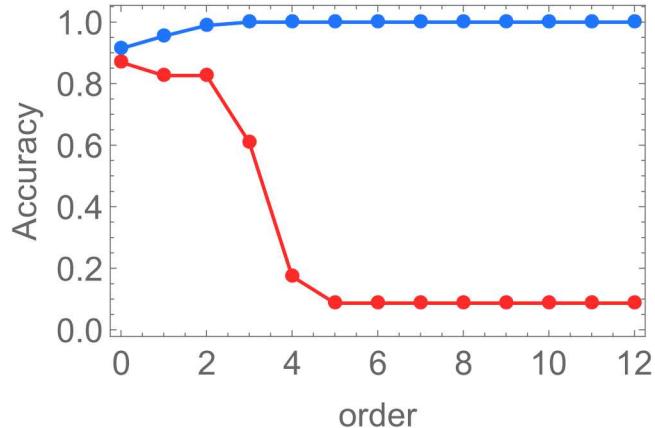


There are interesting remaining problems in advancing the state of the art in this field.



We don't care about compression, so we can modify the algorithms.

Unbalanced Training Sets





Compression algorithms are slow.



“This is a question.”



$\{(' '), ('a'), ('q'), ('.'), ('T'), ('e'), ('h'), ('i'), ('io'), ('is'), ('n'), ('s'), ('st'), ('u')\}$

“This is another question.”



$\{(' '), ('a'), ('q'), ('T'), ('e'), ('h'), ('he'), ('i'), ('io'), ('is'), ('n'), ('n.'), ('o'), ('r'), ('s'), ('st'), ('t'), ('u')\}$

Streaming LZJD is faster when computing sliding windows.



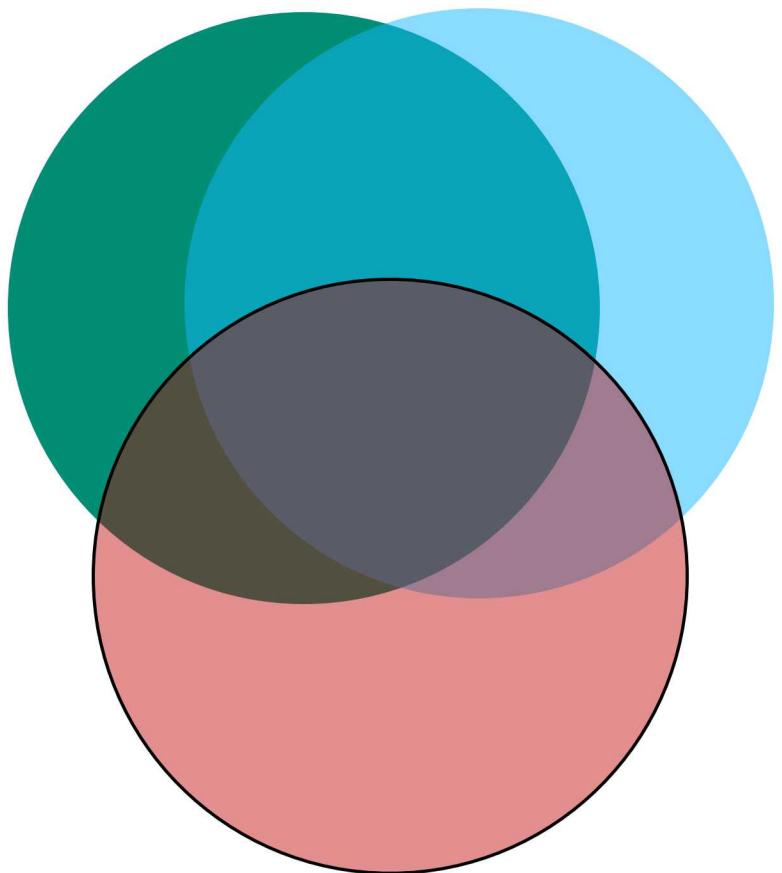
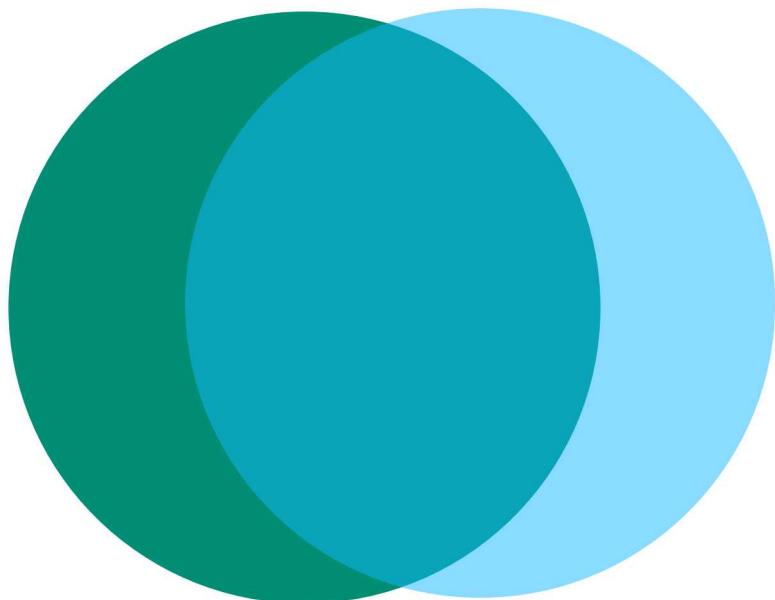
“This is another question.”



```
[('T'), ('h'), ('i'), ('s'), (' '), ('i'), ('is'),  
(' '), ('a'), ('n'), ('o'), ('t'), ('h'),  
('he'), ('r'), (' '), ('q'), ('u'), ('e'),  
('s'), ('st'), ('i'), ('io'), ('n')]
```



You don't have to specify the features, but you can't specify the features.





This is a question.

The diagram illustrates the relationship between natural language and its representation in HTML. At the top, the text "This is a question." is displayed. A teal arrow points from the word "question." to the word "question" in the HTML code below. At the bottom, the HTML code "<p>This is HTML</p>" is shown, enclosed in a teal box with a downward-pointing arrow. This visualizes how a single question in natural language is represented by a single HTML paragraph element.

```
<p>This is HTML</p>
```



Travis Bauer (tlbauer@sandia.gov)

Publications

- McNamara, L., Bauer, T., Haass, M., and Matzen, L., "Silver Ticket? Exploring Entropy Metrics for Visual Analytics." BELIV 2016, Oct 24, 2016, Baltimore.
- Brounstein, Tom Rego, et al. *Stylometric and Temporal Techniques for Social Media Account Resolution*. No. SAND2017-2965C. Sandia National Lab.(SNL-NM), Albuquerque, NM (United States), 2017.
- Ting, C., Fisher, A., Bauer, T., "Compression based algorithms for deception detection." SocInfo2017, Sep 13-15, 2017, Oxford, UK
- Ting, C., Field, R., Fisher, A., Bauer, T., "Compression Analytics for Classification and Anomaly Detection within Network Communication." IEEE Transactions on Information Forensics & Security Volume 14, Issue 5, May 2019
- Ting, C., Field, R., Quach, T., Bauer, T., "Generalized Boundary Detection in Streaming Data Using Compression-Based Analytics. ICASSP 2019, accepted