

HELMHOLTZ  
MUNICH



## What's in a Graph?

Bastian Rieck (@Pseudomanifold)

# Invocation

*What's in a name? That which we call a rose  
By any other name would smell as sweet*



*(Romeo and Juliet, Act II, Scene II)*

# Invocation

*What's in a name? That which we call a rose  
By any other name would smell as sweet*



*(Romeo and Juliet, Act II, Scene II)*

*What's in a graph? That which we call our data  
By any other means would train as well*



# What is a graph?

A graph is a tuple  $(V, E)$ , consisting of a set of vertices  $V$  and a set of edges  $E$ , consisting of subsets of paired vertices.

# What is a graph?

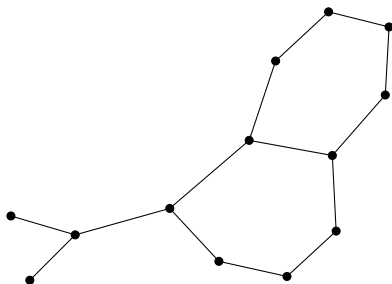
A graph is a tuple  $(V, E)$ , consisting of a set of vertices  $V$  and a set of edges  $E$ , consisting of subsets of paired vertices.

If  $E$  is *ordered* instead, the graph is called *directed*.

# What is a graph?

A graph is a tuple  $(V, E)$ , consisting of a set of vertices  $V$  and a set of edges  $E$ , consisting of subsets of paired vertices.

If  $E$  is *ordered* instead, the graph is called *directed*.



# Alternative views on graphs

A graph is a *triangulation* of a manifold.

# Alternative views on graphs

A graph is a *triangulation* of a manifold.

A graph is a 1-dimensional simplicial complex.



# Alternative views on graphs

A graph is a *triangulation* of a manifold.

A graph is a 1-dimensional simplicial complex.

A graph is a metric space.

# Alternative views on graphs

A graph is a *triangulation* of a manifold.

A graph is a 1-dimensional simplicial complex.

A graph is a metric space.

A graph is a set system.

# Where do graphs come from?

A collection of different attitudes



---

**Alignment**    Belief

---



# Where do graphs come from?

A collection of different attitudes

---

**Alignment**    Belief

---

**Lawful**        Graphs occur only in graph theory.



# Where do graphs come from?

A collection of different attitudes

---

**Alignment**    Belief

---

**Lawful**        Graphs occur only in graph theory.

**Neutral**        Graphs can arise from other data modalities.



# Where do graphs come from?

A collection of different attitudes

---

<b>Alignment</b>	Belief
------------------	--------

---

<b>Lawful</b>	Graphs occur only in graph theory.
---------------	------------------------------------

<b>Neutral</b>	Graphs can arise from other data modalities.
----------------	--

<b>Chaotic</b>	Everything is a graph.
----------------	------------------------

---



# Where do graphs come from?

A collection of different attitudes

---

<b>Alignment</b>	Belief
<b>Lawful</b>	Graphs occur only in graph theory.
<b>Neutral</b>	Graphs can arise from other data modalities.
<b>Chaotic</b>	Everything is a graph.

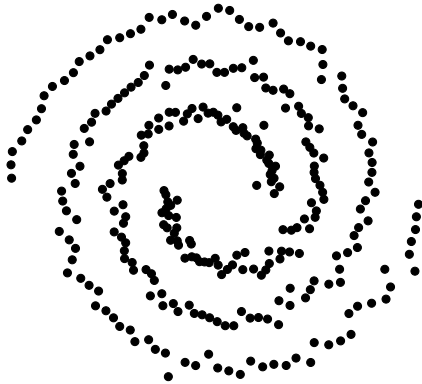
---

*Most graph theorists will agree that among the vast number of graphs that exist there are only a few thousand that can be considered really interesting.*

*(<https://houseofgraphs.org>)*

# Obtaining graphs from other modalities

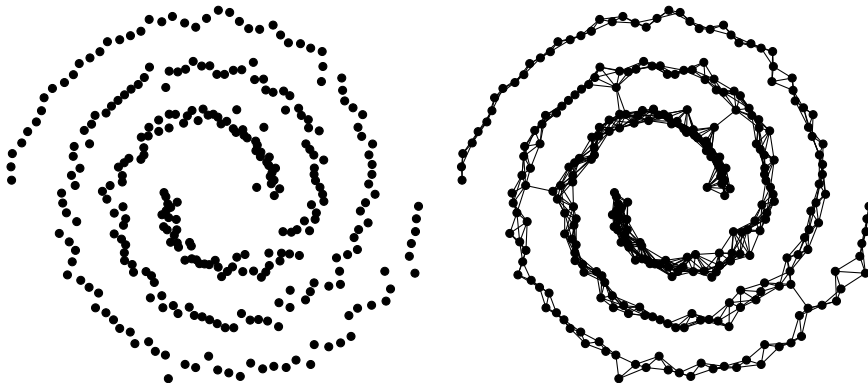
Point clouds





# Obtaining graphs from other modalities

Point clouds

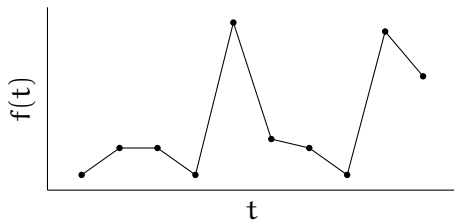


**Rips graph at scale  $\epsilon$**

$$\mathcal{R}_\epsilon := (X, E) \text{ with } E := \{x, y \in X \mid d(x, y) \leq \epsilon\}$$

# Obtaining graphs from other modalities

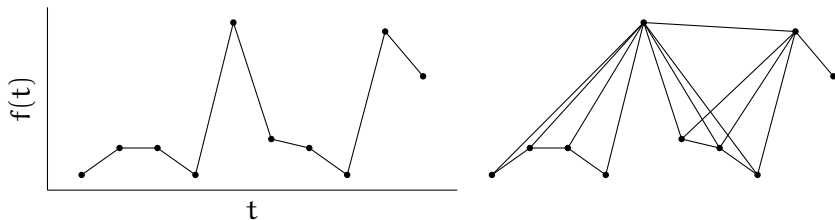
Time series



<sup>1</sup>L. Lacasa, B. Luque, F. Ballesteros, J. Luque and J. C. Nuño, 'From time series to complex networks: The visibility graph', *Proceedings of the National Academy of Sciences* 105.13, 2008, pp. 4972–4975.

# Obtaining graphs from other modalities

Time series



## Visibility graph<sup>1</sup>

Connect observations  $(t_i, f_i)$  and  $(t_{i+1}, f_{i+1})$  if no other observations occur along their linear interpolation.

<sup>1</sup>L. Lacasa, B. Luque, F. Ballesteros, J. Luque and J. C. Nuño, 'From time series to complex networks: The visibility graph', *Proceedings of the National Academy of Sciences* 105.13, 2008, pp. 4972–4975.

# How do these graphs *differ* from other graphs?

Some graphs arise from some prior *geometry*.

# How do these graphs *differ* from other graphs?

Some graphs arise from some prior *geometry*.

Attributes in such graphs can reflect said geometry.

# How do these graphs *differ* from other graphs?

Some graphs arise from some prior *geometry*.

Attributes in such graphs can reflect said geometry.

When studying the graph, we are actually studying its geometry.

# How do these graphs *differ* from other graphs?

Some graphs arise from some prior *geometry*.

Attributes in such graphs can reflect said geometry.

When studying the graph, we are actually studying its geometry.

*Some are born geometrical,  
some achieve geometry,  
and some have geometry thrust upon them.*



We are not equipped to ask what it means to study a specific graph.





We are not equipped to ask what it means to study a specific graph.  
We need to develop a (better) language for describing graph data.



We are not equipped to ask what it means to study a specific graph.

We need to develop a (better) language for describing graph data.

By treating all graphs the same, we are making a mistake.

Does Topology Help in Graph Learning?

# Hypothesis

Graphs are topological objects.

# Hypothesis

Graphs are topological objects.

Graph learning algorithms are not necessarily aware of topological features.

# Hypothesis

Graphs are topological objects.

Graph learning algorithms are not necessarily aware of topological features.

Thus, making models aware of such features is bound to improve predictive performance.

# Making Weisfeiler–Le(h)man persistent

**B. Rieck\***, C. Bock\* and K. Borgwardt, 'A Persistent Weisfeiler–Lehman Procedure for Graph Classification', *Proceedings of the 36th International Conference on Machine Learning (ICML)*, 2019, pp. 5448–5458

Make Weisfeiler–Le(h)man aware of connected components and cycles.

# Making Weisfeiler–Le(h)man persistent

**B. Rieck\***, C. Bock\* and K. Borgwardt, 'A Persistent Weisfeiler–Lehman Procedure for Graph Classification', *Proceedings of the 36th International Conference on Machine Learning (ICML)*, 2019, pp. 5448–5458

Make Weisfeiler–Le(h)man aware of connected components and cycles.

Show that the known algorithm is actually one a specific instance of a larger scheme.



# Making Weisfeiler–Le(h)man persistent

**B. Rieck\***, C. Bock\* and K. Borgwardt, 'A Persistent Weisfeiler–Lehman Procedure for Graph Classification', *Proceedings of the 36th International Conference on Machine Learning (ICML)*, 2019, pp. 5448–5458

Make Weisfeiler–Le(h)man aware of connected components and cycles.

Show that the known algorithm is actually one a specific instance of a larger scheme.

Build intuition based on quick, iterative development.

# Making Weisfeiler–Le(h)man persistent

B. Rieck\*, C. Bock\* and K. Borgwardt, 'A Persistent Weisfeiler–Lehman Procedure for Graph Classification', *Proceedings of the 36th International Conference on Machine Learning (ICML)*, 2019, pp. 5448–5458

Make Weisfeiler–Le(h)man aware of connected components and cycles.

Show that the known algorithm is actually one a specific instance of a larger scheme.

Build intuition based on quick, iterative development.

<i>Method</i>	<i>Data set</i>	
	MUTAG	PTC-MM
VH	85.96 ± 0.27	66.96 ± 0.51
EH	85.69 ± 0.46	61.61 ± 0.00
WL	87.26 ± 1.42	67.28 ± 0.97
P-WL	86.10 ± 1.37	<b>68.40 ± 1.17</b>
P-WL-C	<b>90.51 ± 1.34</b>	<b>68.57 ± 1.76</b>

# Making graph neural networks topology-aware

M. Horn\*, E. De Brouwer\*, M. Moor, Y. Moreau, **B. Rieck**<sup>†</sup> and K. Borgwardt<sup>†</sup>, 'Topological Graph Neural Networks', *International Conference on Learning Representations (ICLR)*, 2022, arXiv: 2102.07835 [cs.LG]

Providing a new layer for use in graph neural networks.

# Making graph neural networks topology-aware

M. Horn\*, E. De Brouwer\*, M. Moor, Y. Moreau, **B. Rieck**<sup>†</sup> and K. Borgwardt<sup>†</sup>, 'Topological Graph Neural Networks', *International Conference on Learning Representations (ICLR)*, 2022, arXiv: 2102.07835 [cs.LG]

Providing a new layer for use in graph neural networks.

Create new synthetic data sets with pronounced topological features.

# Making graph neural networks topology-aware

M. Horn\*, E. De Brouwer\*, M. Moor, Y. Moreau, **B. Rieck**<sup>†</sup> and K. Borgwardt<sup>†</sup>, 'Topological Graph Neural Networks', *International Conference on Learning Representations (ICLR)*, 2022, arXiv: 2102.07835 [cs.LG]

Providing a new layer for use in graph neural networks.

Create new synthetic data sets with pronounced topological features.

Getting high predictive performance on these data sets.



# Making graph neural networks topology-aware

M. Horn\*, E. De Brouwer\*, M. Moor, Y. Moreau, **B. Rieck**<sup>†</sup> and K. Borgwardt<sup>†</sup>, 'Topological Graph Neural Networks', *International Conference on Learning Representations (ICLR)*, 2022, arXiv: 2102.07835 [cs.LG]

Providing a new layer for use in graph neural networks.

Create new synthetic data sets with pronounced topological features.

Getting high predictive performance on these data sets.

Getting mediocre performance on other benchmark data sets.



# Making graph neural networks topology-aware

M. Horn\*, E. De Brouwer\*, M. Moor, Y. Moreau, **B. Rieck**<sup>†</sup> and K. Borgwardt<sup>†</sup>, 'Topological Graph Neural Networks', *International Conference on Learning Representations (ICLR)*, 2022, arXiv: 2102.07835 [cs.LG]

Providing a new layer for use in graph neural networks.

Create new synthetic data sets with pronounced topological features.

Getting high predictive performance on these data sets.

Getting mediocre performance on other benchmark data sets.

Getting better performance if node features are randomised.



*On data sets with pronounced topological structures, we found that our method helps GNNs obtain substantial gains in predictive performance.*

Does Topology Help in Graph Learning?

*It depends.*



# Lessons learned

Make your baseline as strong as possible.

# Lessons learned

Make your baseline as strong as possible.

Iterate quickly to check your hypotheses.

# Lessons learned

Make your baseline as strong as possible.

Iterate quickly to check your hypotheses.

Ablation studies are crucial.

# Lessons learned

Make your baseline as strong as possible.

Iterate quickly to check your hypotheses.

Ablation studies are crucial.

Use *repeated nested cross-validation* for small data sets.

L. O'Bray\*, **B. Rieck\*** and K. Borgwardt, 'Filtration Curves for Graph Representation', *Proceedings of the 27th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining (KDD)*, 2021, pp. 1267–1275:

*For example, [our] previous work had results as high as 80% on IMDB-BINARY when considering just a single run of 10-fold cross validation. However, the results were not reflective of performance when repeated 10 times, which reduced [performance] to around 73%.*

# Troubling Trends

## Concerning baselines (I)

K. Borgwardt, E. Ghisu, F. Llinares-López, L. O'Bray and **B. Rieck**, 'Graph Kernels: State-of-the-Art and Future Challenges', *Foundations and Trends® in Machine Learning* 13.5–6, 2020, pp. 531–712, arXiv: 2011.03854 [cs.LG]

We analysed the performance of *graph kernels* on benchmark data sets.

## Concerning baselines (I)

K. Borgwardt, E. Ghisu, F. Llinares-López, L. O'Bray and **B. Rieck**, 'Graph Kernels: State-of-the-Art and Future Challenges', *Foundations and Trends® in Machine Learning* 13.5–6, 2020, pp. 531–712, arXiv: 2011.03854 [cs.LG]

We analysed the performance of *graph kernels* on benchmark data sets.

Histogram kernels turn out to be surprisingly effective.



## Concerning baselines (I)

K. Borgwardt, E. Ghisu, F. Llinares-López, L. O'Bray and **B. Rieck**, 'Graph Kernels: State-of-the-Art and Future Challenges', *Foundations and Trends® in Machine Learning* 13.5–6, 2020, pp. 531–712, arXiv: 2011.03854 [cs.LG]

We analysed the performance of *graph kernels* on benchmark data sets.

Histogram kernels turn out to be surprisingly effective.

Not using any deeper insights into graph structure here.



<i>Data set</i>	VH	EH
AIDS	99.70 ± 0.00	99.28 ± 0.07
DD	68.68 ± 3.04	78.52 ± 0.34
IMDB-BINARY	50.58 ± 0.20	73.46 ± 0.60
MUTAG	85.98 ± 0.40	85.14 ± 1.03
Mutagenicity	67.01 ± 0.83	49.13 ± 1.75
NCI1	64.66 ± 0.53	51.71 ± 1.45
NCI109	63.24 ± 0.54	51.45 ± 2.06
REDDIT-BINARY	50.03 ± 2.24	78.94 ± 0.60
SYNTHETICnew	62.30 ± 0.55	71.20 ± 1.69



## Concerning baselines (II)

C. Cai and Y. Wang, 'A simple yet effective baseline for non-attributed graph classification', 2018, arXiv: 1811.03508 [cs.LG]

Use a *local degree profile* to classify graphs, then train an SVM on the resulting features.

## Concerning baselines (II)

C. Cai and Y. Wang, 'A simple yet effective baseline for non-attributed graph classification', 2018, arXiv: 1811.03508 [cs.LG]

Use a *local degree profile* to classify graphs, then train an SVM on the resulting features.

The features turn out to be surprisingly effective.



## Concerning baselines (II)

C. Cai and Y. Wang, 'A simple yet effective baseline for non-attributed graph classification', 2018, arXiv: 1811.03508 [cs.LG]

Use a *local degree profile* to classify graphs, then train an SVM on the resulting features.

The features turn out to be surprisingly effective.



Not using any deeper insights into graph structure here.



*Most graph kernels aim to capture graph topology and graph similarity in the hope of improving classification accuracy. Our experiment[s] suggests that this is not yet well-reflected on current benchmark datasets for non-attribute[d] graphs.*

## Concerning baselines (II)

C. Cai and Y. Wang, 'A simple yet effective baseline for non-attributed graph classification', 2018, arXiv: 1811.03508 [cs.LG]

Use a *local degree profile* to classify graphs, then train an SVM on the resulting features.

The features turn out to be surprisingly effective.



Not using any deeper insights into graph structure here.



*Most graph kernels aim to capture graph topology and graph similarity in the hope of improving classification accuracy. Our experiment[s] suggests that this is not yet well-reflected on current benchmark datasets for non-attribute[d] graphs.*

*In general, while not addressed in this paper, we note that understanding the power and limitation of various graph representations, [...] is crucial [...] and remains largely open.*



# Is this graph the *right* graph?

J. Gasteiger, S. Weißberger and S. Günnemann, 'Diffusion Improves Graph Learning', *Advances in Neural Information Processing Systems*, vol. 32, Curran Associates, Inc., 2019:

*Edges in real graphs are often noisy or defined using an arbitrary threshold, so we can clearly improve upon this approach.*

J. Topping, F. Di Giovanni, B. P. Chamberlain, X. Dong and M. M. Bronstein, 'Understanding over-squashing and bottlenecks on graphs via curvature', *International Conference on Learning Representations*, 2022:

*More recently, there is a trend to decouple the input graph from the graph used for information propagation.*

## Question

Is the graph structure not necessary and we could equally well solve everything with a properly-regularised transformer-like architecture?

## A Closer Look At Our Data

## Our data sets are not necessarily *representative*

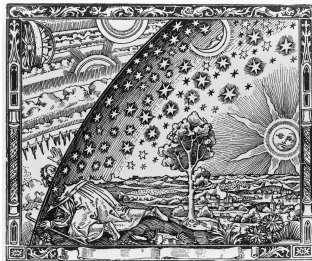
J. Palowitch, A. Tsitsulin, B. Mayer and B. Perozzi, 'GraphWorld: Fake Graphs Bring Real Insights for GNNs', *Proceedings of the 28th ACM SIGKDD Conference on Knowledge Discovery & Data Mining (KDD)*, 2022, pp. 3691–3701

*Oh, God, I could be bounded in a nutshell and count myself a king of infinite space, were it not that I have bad dreams.*  
(Hamlet, Act II, Scene II)

# Our data sets are not necessarily *representative*

J. Palowitch, A. Tsitsulin, B. Mayer and B. Perozzi, 'GraphWorld: Fake Graphs Bring Real Insights for GNNs', *Proceedings of the 28th ACM SIGKDD Conference on Knowledge Discovery & Data Mining (KDD)*, 2022, pp. 3691–3701

*Oh, God, I could be bounded in a nutshell and count myself a king of infinite space, were it not that I have bad dreams.*  
(Hamlet, Act II, Scene II)



*Our first finding is that standard benchmark graphs [...] cover only a small region of this graph space that GraphWorld is able to cover via synthetic graph generation.*



# Where do our data sets come from?

*And some things that should not have been forgotten were lost. History became legend.  
Legend became myth.* *(The Lord of the Rings)*

# Where do our data sets come from?

*And some things that should not have been forgotten were lost. History became legend.  
Legend became myth.* *(The Lord of the Rings)*

Typically, no *provenance* information of (benchmark) data sets.

# Where do our data sets come from?

*And some things that should not have been forgotten were lost. History became legend.  
Legend became myth.* *(The Lord of the Rings)*

Typically, no *provenance* information of (benchmark) data sets.

Typically, no *version* information.

# Where do our data sets come from?

*And some things that should not have been forgotten were lost. History became legend.  
Legend became myth.* *(The Lord of the Rings)*

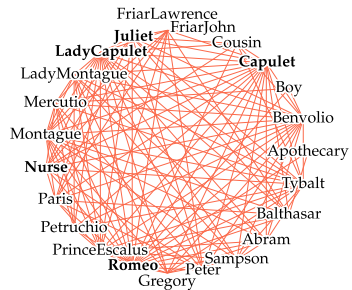
Typically, no *provenance* information of (benchmark) data sets.

Typically, no *version* information.

Often, no pre-defined splits.

# Why provenance information is important

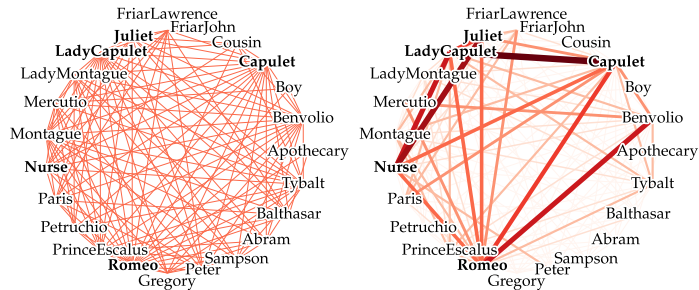
C. Coupette, J. Vreeken and **B. Rieck**, 'All the World's a (Hyper)Graph: A Data Drama', 2022, arXiv: 2206.08225 [cs.LG], URL: <https://hyperbard.net>



Three *valid* co-occurrence networks of characters in Shakespeare's Romeo and Juliet. Characters in Act III, Scene V are highlighted.

# Why provenance information is important

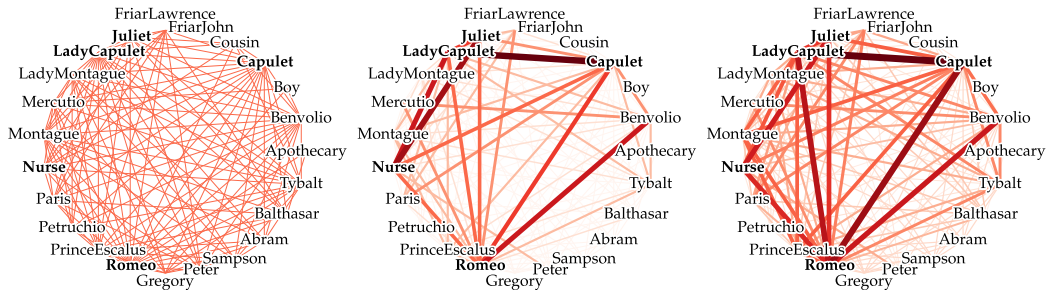
C. Coupette, J. Vreeken and **B. Rieck**, 'All the World's a (Hyper)Graph: A Data Drama', 2022, arXiv: 2206.08225 [cs.LG], URL: <https://hyperbard.net>



Three *valid* co-occurrence networks of characters in Shakespeare's Romeo and Juliet. Characters in Act III, Scene V are highlighted.

# Why provenance information is important

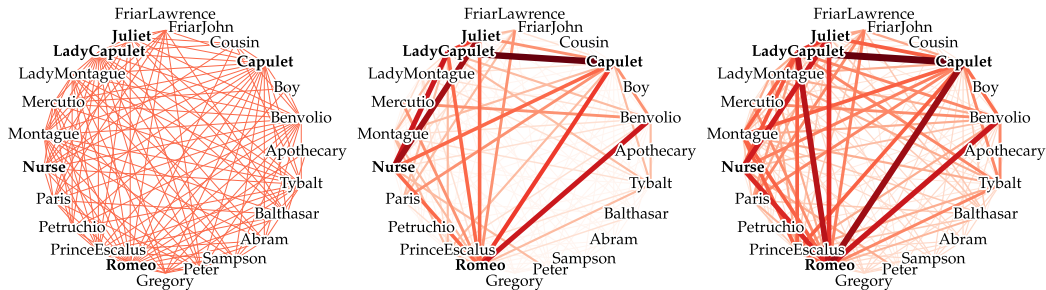
C. Coupette, J. Vreeken and **B. Rieck**, 'All the World's a (Hyper)Graph: A Data Drama', 2022, arXiv: 2206.08225 [cs.LG], URL: <https://hyperbard.net>



Three *valid* co-occurrence networks of characters in Shakespeare's *Romeo and Juliet*. Characters in Act III, Scene V are highlighted.

# Why provenance information is important

C. Coupette, J. Vreeken and **B. Rieck**, 'All the World's a (Hyper)Graph: A Data Drama', 2022, arXiv: 2206.08225 [cs.LG], URL: <https://hyperbard.net>



Three *valid* co-occurrence networks of characters in Shakespeare's *Romeo and Juliet*. Characters in Act III, Scene V are highlighted.

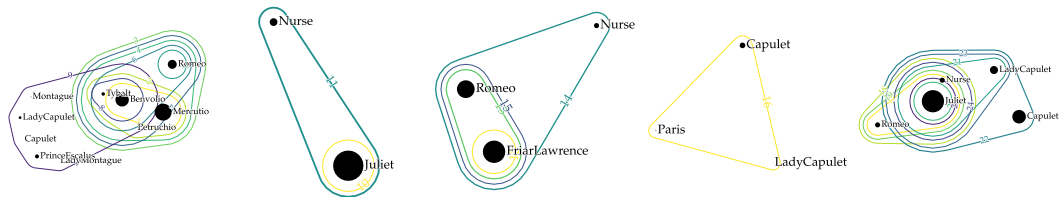
## Observation

Romeo and the Capulets almost never interact directly; our *modelling decision* introduces new information!



# Alternative models

## Hypergraphs



Modelling individual scenes of Act III of Romeo and Juliet provides a better overview of how all characters interact on stage.

# Where do our data sets come from?

pytorch-geometric

**Disclaimer:** The PyG team is doing an excellent job! This is *not* to be construed as a criticism of their work. We, as a community, should do more to support such endeavours!

```
$ rg -t py "^\\s*url =" \
  | grep -Eo "(http|https)://[a-zA-Z0-9./?=_%:-]*" \
  | sort \
  | uniq -c
```

## Hosts for data sets

<i>Host</i>	<i>Count</i>
ucl.ac.uk	1
is.tue.mpg.de	3
docs.google.com	7
drive.google.com	7
dropbox.com	7

# Where do our data sets come from?

pytorch-geometric, continued

Some additional data sets are also accessed via GitHub repositories, owned by private individuals or organisations:

abojchevski, flyingdoog, gasteigerjo, graphdml-uiuc-jlu, nd7141, INK-USC, kimiyoung, klicperajo, pmernyei, samihaija, shchur, steveazzolin, villmow, yandex-research, Yannick-S

# Where do our data sets come from?

pytorch-geometric, continued

Some additional data sets are also accessed via GitHub repositories, owned by private individuals or organisations:

abojchevski, flyingdoog, gasteigerjo, graphdm1-uiuc-jlu, nd7141, INK-USC, kimiyoung, klicperajo, pmernyei, samihaija, shchur, steveazzolin, villmow, yandex-research, Yannick-S

## Some issues

At present, no versioning or fingerprinting mechanism appears to be in place.

# Where do our data sets come from?

pytorch-geometric, continued

Some additional data sets are also accessed via GitHub repositories, owned by private individuals or organisations:

abojchevski, flyingdoog, gasteigerjo, graphdm1-uiuc-jlu, nd7141, INK-USC, kimiyoung, klicperajo, pmernyei, samihaija, shchur, steveazzolin, villmow, yandex-research, Yannick-S

## Some issues

At present, no versioning or fingerprinting mechanism appears to be in place.

What happens if a repository is deleted?

# Where do our data sets come from?

pytorch-geometric, continued

Some additional data sets are also accessed via GitHub repositories, owned by private individuals or organisations:

abojchevski, flyingdoog, gasteigerjo, graphdm1-uiuc-jlu, nd7141, INK-USC, kimiyoung, klicperajo, pmernyei, samihaija, shchur, steveazzolin, villmow, yandex-research, Yannick-S

## Some issues

At present, no versioning or fingerprinting mechanism appears to be in place.

What happens if a repository is deleted?

What happens if files are being changed?

# Where do our data sets come from?

pytorch-geometric, continued

Some additional data sets are also accessed via GitHub repositories, owned by private individuals or organisations:

abojchevski, flyingdoog, gasteigerjo, graphdm1-uiuc-jlu, nd7141, INK-USC, kimiyoung, klicperajo, pmernyei, samihaija, shchur, steveazzolin, villmow, yandex-research, Yannick-S

## Some issues

At present, no versioning or fingerprinting mechanism appears to be in place.

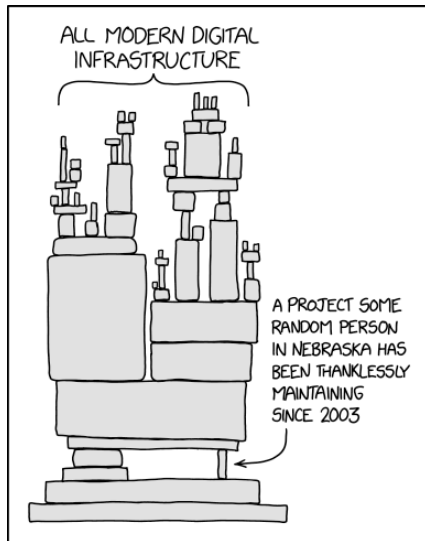
What happens if a repository is deleted?

What happens if files are being changed?

Is everyone training on the same data?

# Is this the state of our data sets?

<https://xkcd.com/2347/>





# Is this where we are heading?



B. Haibe-Kains et al., 'Transparency and reproducibility in artificial intelligence', *Nature* 586.7829, 2020, E14–E16

## Matters arising

# Transparency and reproducibility in artificial intelligence

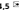
<https://doi.org/10.1038/s41586-020-2766-y>

Received: 1 February 2020

Accepted: 10 August 2020

Published online: 14 October 2020

 Check for updates

Benjamin Haibe-Kains<sup>12,3,4,5</sup> , George Alexandru Adam<sup>3,5</sup>, Ahmed Hosny<sup>6,7</sup>, Farnoosh Khodakarami<sup>1,2</sup>, Massive Analysis Quality Control (MAQC) Society Board of Directors\*, Levi Waldron<sup>8</sup>, Bo Wang<sup>2,3,5,9,10</sup>, Chris McIntosh<sup>2,5,9</sup>, Anna Goldenberg<sup>3,5,11,12</sup>, Anshul Kundaje<sup>13,14</sup>, Casey S. Greene<sup>15,16</sup>, Tamara Broderick<sup>17</sup>, Michael M. Hoffman<sup>1,2,3,5</sup>, Jeffrey T. Leek<sup>18</sup>, Keegan Korthauer<sup>19,20</sup>, Wolfgang Huber<sup>21</sup>, Alvis Brazma<sup>22</sup>, Joelle Pineau<sup>23,24</sup>, Robert Tibshirani<sup>25,26</sup>, Trevor Hastie<sup>25,26</sup>, John P. A. Ioannidis<sup>25,26,27,28,29</sup>, John Quackenbush<sup>30,31,32</sup> & Hugo J. W. L. Aerts<sup>6,7,33,34</sup>

arising from S. M. McKinney et al. *Nature* <https://doi.org/10.1038/s41586-019-1799-6> (2020)

Improving the Field

# What can we do about that?

Keeping track of the provenance of data sets.

<sup>2</sup>T. Gebru et al., 'Datasheets for Datasets', *Communications of the ACM* 64.12, 2021, pp. 86–92.

# What can we do about that?

Keeping track of the provenance of data sets.

Fingerprinting of data sets (SHA-256).

<sup>2</sup>T. Gebru et al., 'Datasheets for Datasets', *Communications of the ACM* 64.12, 2021, pp. 86–92.

# What can we do about that?

Keeping track of the provenance of data sets.

Fingerprinting of data sets (SHA-256).

Versioning data sets.

<sup>2</sup>T. Gebru et al., 'Datasheets for Datasets', *Communications of the ACM* 64.12, 2021, pp. 86–92.

# What can we do about that?

Keeping track of the provenance of data sets.

Fingerprinting of data sets (SHA-256).

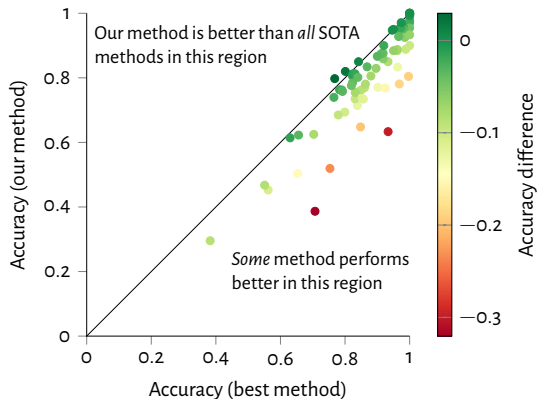
Versioning data sets.

Using *datasheets* to describe graph data sets.<sup>2</sup>

<sup>2</sup>T. Gebru et al., 'Datasheets for Datasets', *Communications of the ACM* 64.12, 2021, pp. 86–92.

# Reporting results

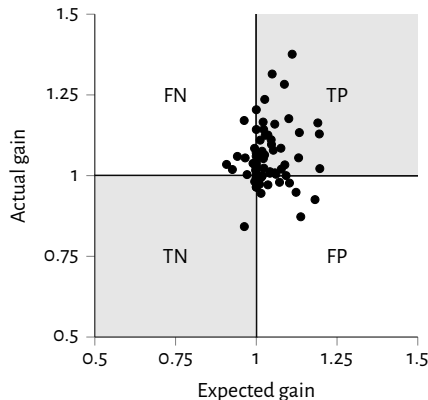
Learning from time series classification in data mining: run comparisons on as many data sets as possible, always comparing with the best available method.



Comparison of the predictive accuracy of our method against the respective state-of-the-art method for the 'UCR Time Series Archive' data sets.

# Reporting results

Learning from time series classification in data mining: calculate 'Texas Sharpshooter' plots, using a well-established algorithm as a baseline.<sup>3</sup>



A 'Texas Sharpshooter' plot, comparing *expected* gains (measured by comparing training performance) of our method with *actual* gains (measured by comparing test performance), relative to 1-DTW-KNN.

<sup>3</sup>G. E. A. P. A. Batista, E. J. Keogh, O. M. Tataw and V. M. A. de Souza, 'CID: an efficient complexity-invariant distance



# Ending on a good note

Our problems are less of a technological and more of a social nature!

# Ending on a good note

Our problems are less of a technological and more of a social nature!  
We can fix them together.

# Ending on a good note

Our problems are less of a technological and more of a social nature!  
We can fix them together.

## Acknowledgements

A lot of people helped shaped the opinions presented in this talk. Thank you to all of them!



Image sources: [xkcd.com](http://xkcd.com); Gerard Girbes Berges (The Noun Project)