# Introduction to Machine Learning for the Digital Humanities

**Klaus-Robert Müller !!et al.!!**

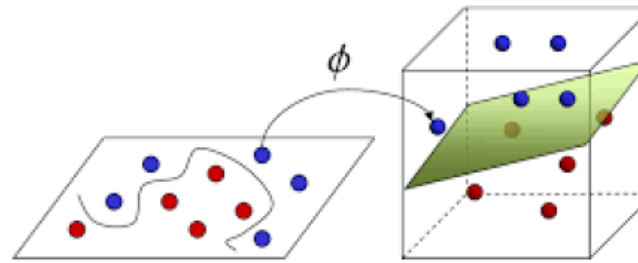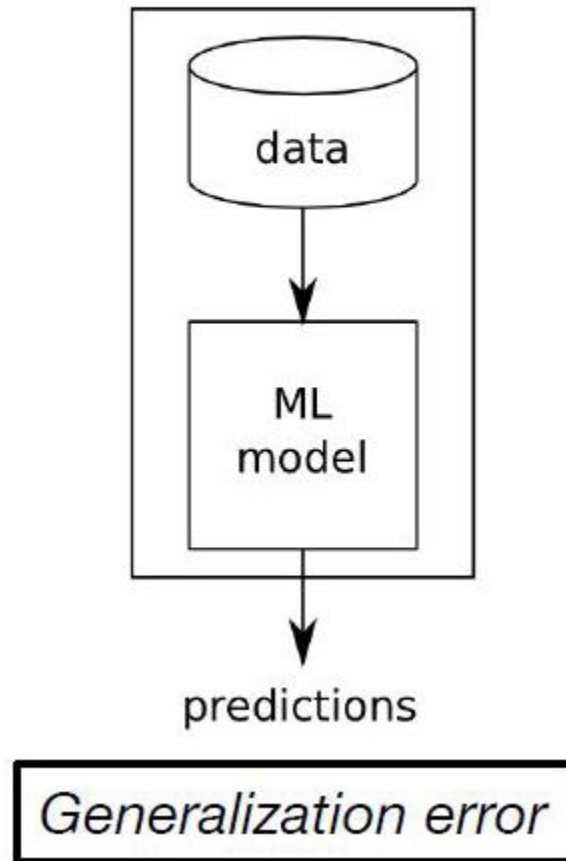# Outline

- A brief introduction to ML concepts

- Remarks on interdisciplinary collaborations of ML with *

- Applications in Neuroscience and Physics

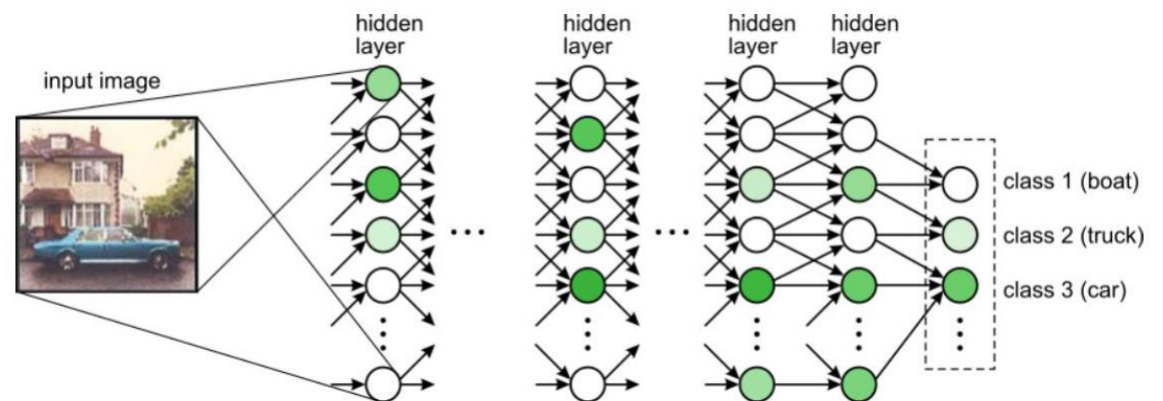- selected recent contributions of ML and for DH

# ML in a nutshell

## Standard ML



data

ML model

predictions

*Generalization error*

## Kernel Methods: SVM etc.



$$\max_{\boldsymbol{\alpha}} \quad W(\boldsymbol{\alpha}) = \sum_{i=1}^{\ell} \alpha_i - \frac{1}{2} \sum_{i,j=1}^{\ell} \alpha_i \alpha_j y_i y_j \, \mathrm{k}(\boldsymbol{x}_i, \boldsymbol{x}_j)$$

subject to $\quad 0 \le \alpha_i \le C, \ i = 1, \ldots, \ell, \ \text{and} \ \sum_{i=1}^{\ell} \alpha_i y_i = 0,$

## Deep Neural Networks



input image

hidden layer | hidden layer | hidden layer | hidden layer

class 1 (boat)
class 2 (truck)
class 3 (car)

# Towards Explaining:
# Machine Learning = black box?

# Explaining single Predictions Pixel-wise



input image

**Forward Propagation**

$$x_j = \mathrm{sigm}\left(\sum_i x_i w_{ij}\right)$$

"It's a rooster"

**Relevance Propagation**
(Bach et al. 2015)

$$R_i = \sum_j R_j \frac{x_i w_{ij}}{\sum_i x_i w_{ij}}$$
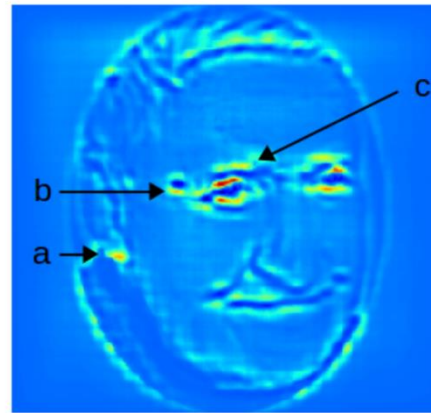
heatmap

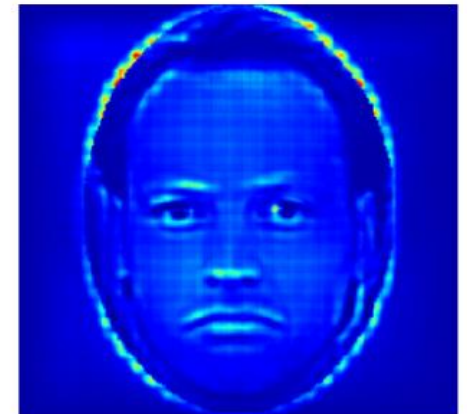**Explaining single decisions is difficult!**

# Applying Explanation in Vision and Text
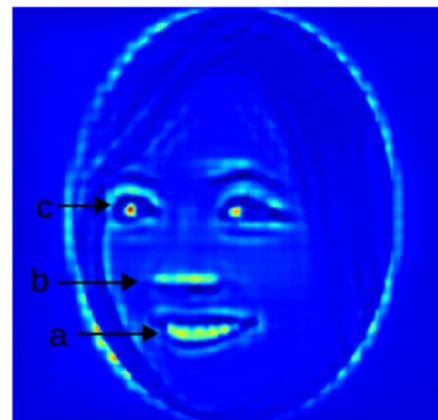
# Application: Faces

What makes
you look old ?

What makes
you look sad ?

What makes
you look attractive ?

# Application: Document Classification

**sci.space**

It is the body's reaction to a strange environment. It appears to be induced partly to physical discomfort and part to mental distress. Some people are more prone to it than others, like some people are more prone to get sick on a roller coaster ride than others. The mental part is usually induced by a lack of clear indication of which way is up or down, ie: the Shuttle is normally oriented with its cargo bay pointed towards Earth, so the Earth (or ground) is "above" the head of the astronauts. About 50% of the astronauts experience some form of motion sickness, and NASA has done numerous tests in space to try to see how to keep the number of occurances down.

**rec.motorcycles**

It is the body's reaction to a strange environment. It appears to be induced partly to physical discomfort and part to mental distress. Some people are more prone to it than others, like some people are more prone to get sick on a roller coaster ride than others. The mental part is usually induced by a lack of clear indication of which way is up or down, ie: the Shuttle is normally oriented with its cargo bay pointed towards Earth, so the Earth (or ground) is "above" the head of the astronauts. About 50% of the astronauts experience some form of motion sickness, and NASA has done numerous tests in space to try to see how to keep the number of occurances down.
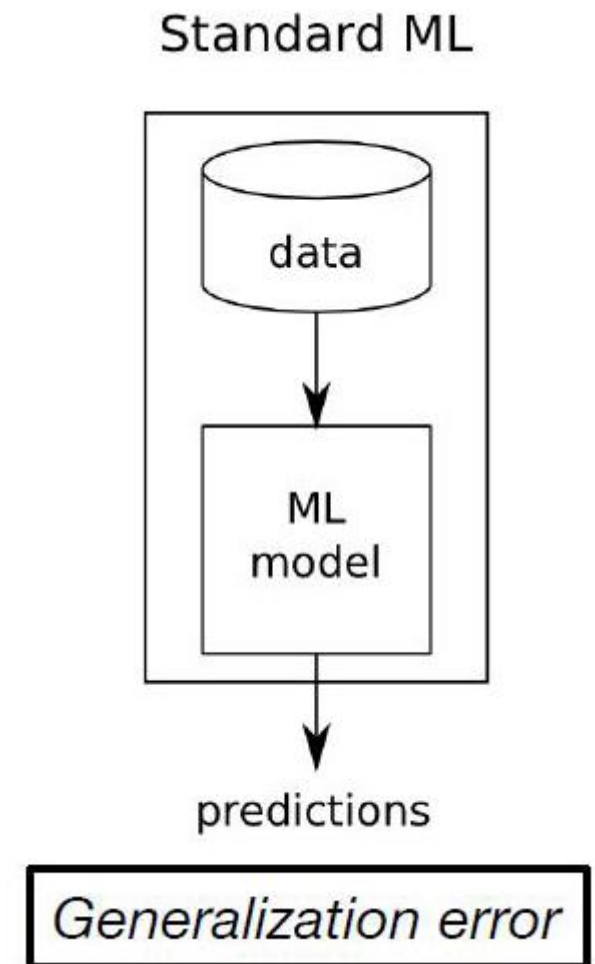
**sci.med**

It is the body's reaction to a strange environment. It appears to be induced partly to physical discomfort and part to mental distress. Some people are more prone to it than others, like some people are more prone to get sick on a roller coaster ride than others. The mental part is usually induced by a lack of clear indication of which way is up or down, ie: the Shuttle is normally oriented with its cargo bay pointed towards Earth, so the Earth (or ground) is "above" the head of the astronauts. About 50% of the astronauts experience some form of motion sickness, and NASA has done numerous tests in space to try to see how to keep the number of occurances down.

# Interlude

# Is the Generalization Error all we need?


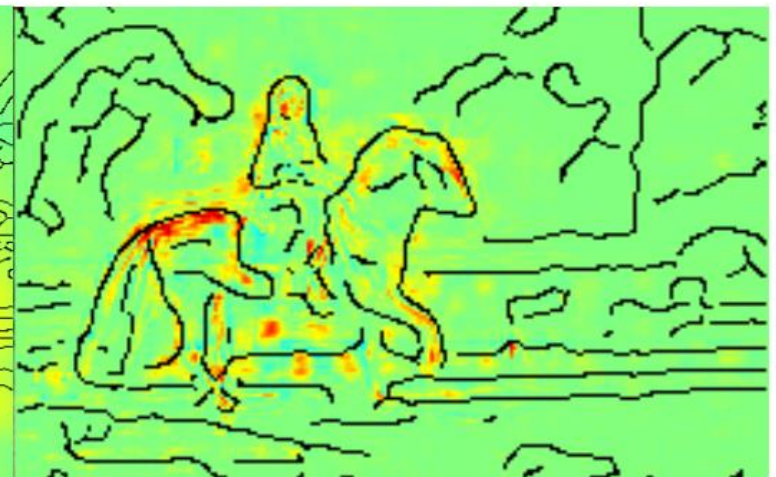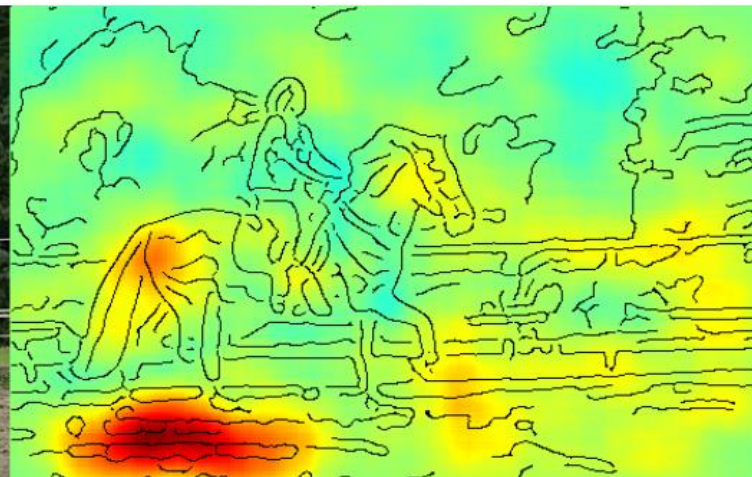Standard ML

# Application: Comparing Classifiers (Lapuschkin et al CVPR 2016)

Test error for various classes:

| | aeroplane | bicycle | bird | boat | bottle | bus | car |
|---|---|---|---|---|---|---|---|
| **Fisher** | 79.08% | 66.44% | 45.90% | 70.88% | 27.64% | 69.67% | 80.96% |
| **DeepNet** | 88.08% | 79.69% | 80.77% | 77.20% | 35.48% | 72.71% | 86.30% |
| | **cat** | **chair** | **cow** | **diningtable** | **dog** | **horse** | **motorbike** |
| **Fisher** | 59.92% | 51.92% | 47.60% | 58.06% | 42.28% | 80.45% | 69.34% |
| **DeepNet** | 81.10% | 51.04% | 61.10% | 64.62% | 76.17% | 81.60% | 79.33% |
| | **person** | **pottedplant** | **sheep** | **sofa** | **train** | **tvmonitor** | **mAP** |
| **Fisher** | 85.10% | 28.62% | 49.58% | 49.31% | 82.71% | 54.33% | 59.99% |
| **DeepNet** | 92.43% | 49.99% | 74.04% | 49.48% | 87.07% | 67.08% | 72.12% |



Image          FV          DNN

# Machine Learning in the Sciences

# Machine Learning in Neuroscience

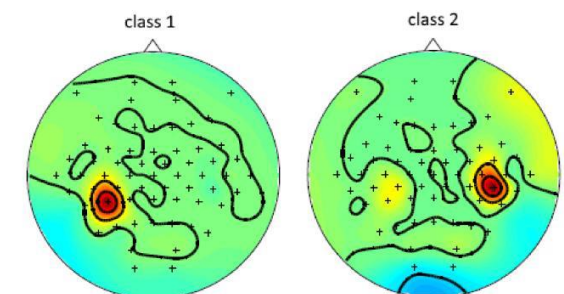# Brain Computer Interfacing: ‚Brain Pong‘



**Berlin Brain Computer Ínterface**

- ML reduces patient training from 300h -> 5min

**Applications**

- help/hope for patients (ALS, stroke…)

- neuroscience

- neurotechnology (video coding, gaming, monitoring driving)

**Leitmotiv: ›let the machines learn‹**

class 1          class 2

# ML4 Quantum Chemistry

# Machine Learning in Chemistry, Physics and Materials

Matthias Rupp, Anatole von Lilienfeld,
Alexandre Tkatchenko, Klaus-Robert Müller

[Rupp et al. Phys Rev Lett 2012, Snyder et al. Phys Rev Lett 2012, Hansen et al. JCTC 2013 and JPCL 2015]

*Ansatz:*

$$\{Z_I, \mathbf{R}_I\} \overset{\mathrm{ML}}{\longmapsto} E$$
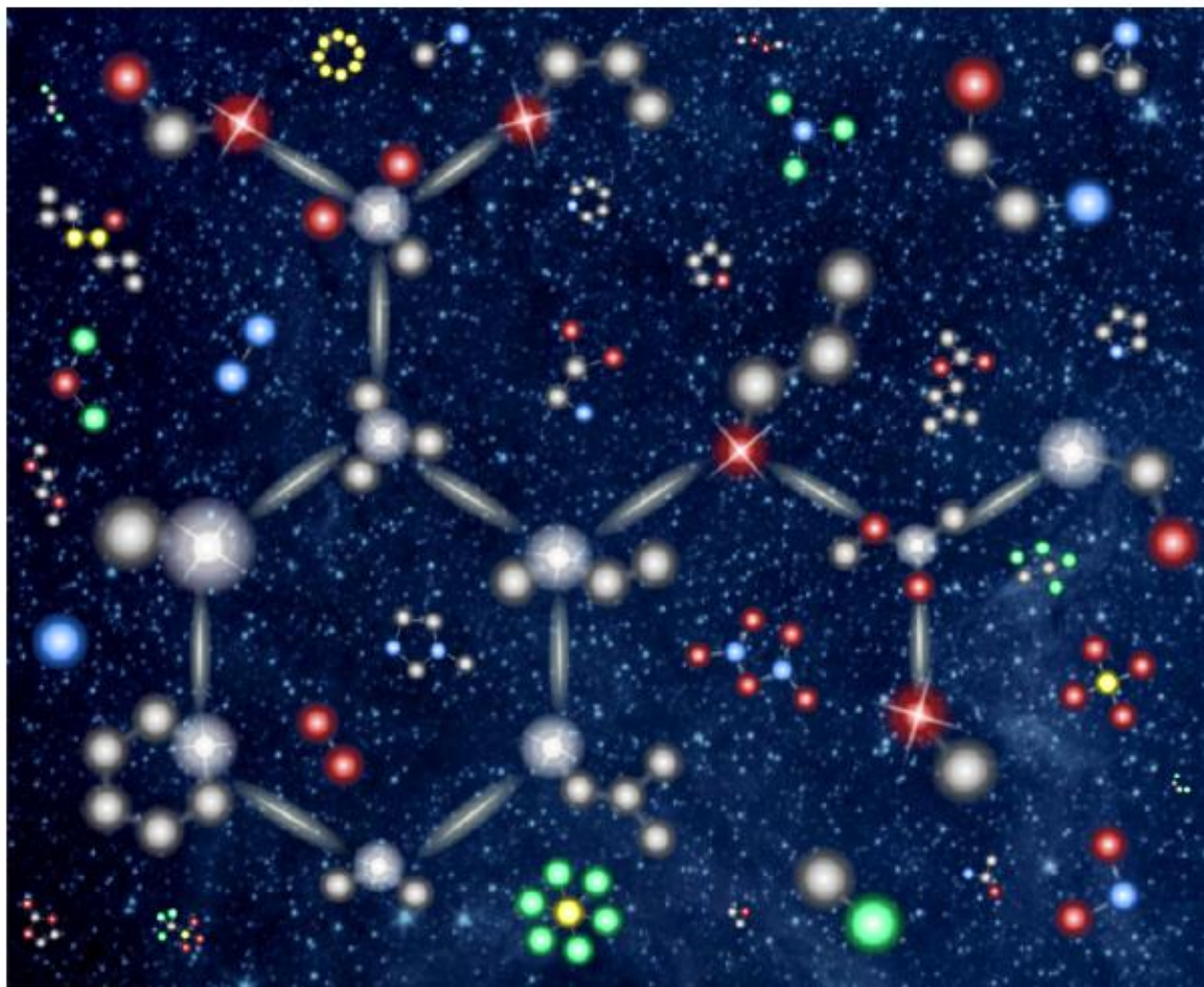
instead of

$$\hat{H}(\{Z_I, \mathbf{R}_I\}) \overset{\Psi}{\longmapsto} E$$
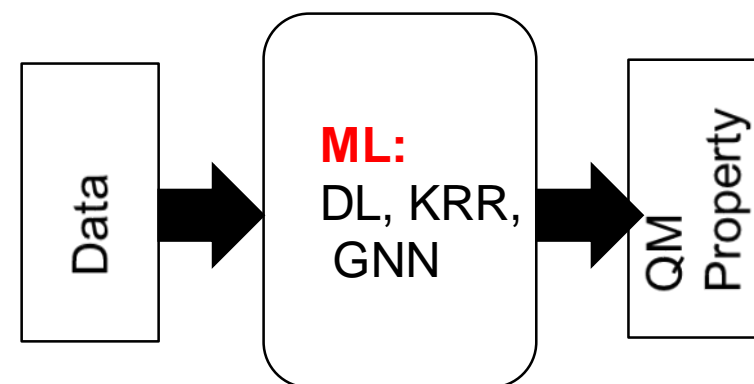
$$\hat{H}\Psi = E\Psi$$



[from von Lilienfeld]

# Navigating Chemical Compound Space



**QC**: Standard **DFT**
3-5h (molecules), 4 months (materials)

**CCSD(T)**: 7days (molecules)

**ML**: 0.1ms

Data → **ML:** DL, KRR, GNN → QM Property
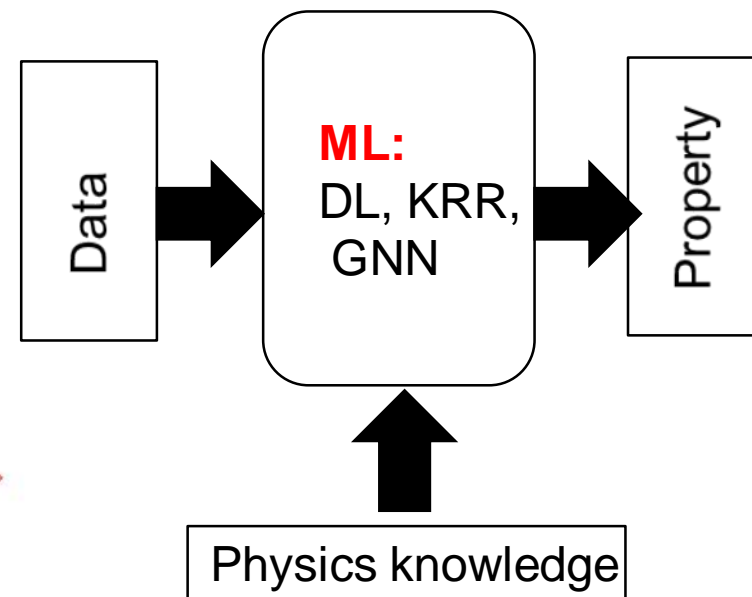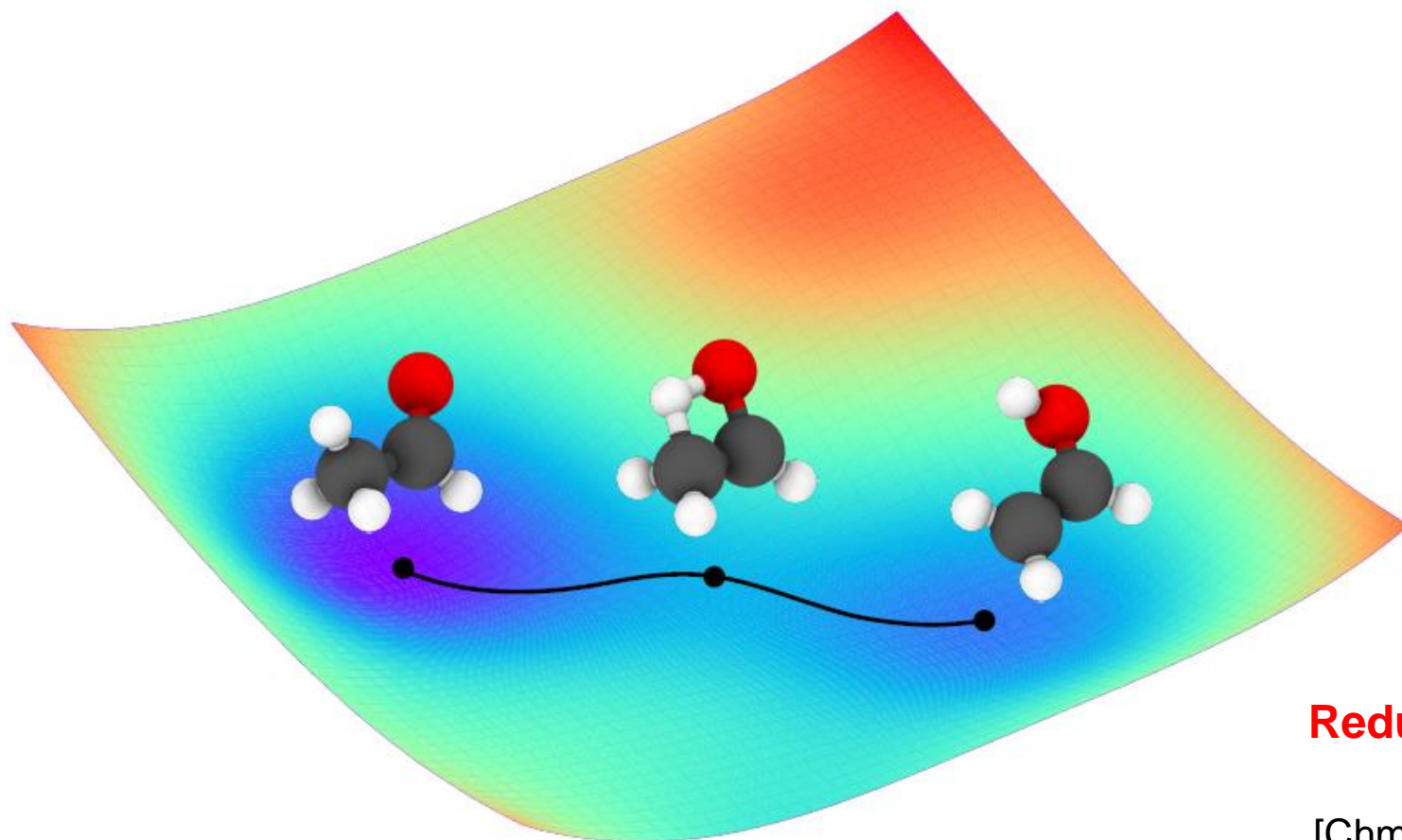
[v Lilienfeld, Tkatchenko, Müller, Nat Rev Chem 2020, Unke et al Chem Rev 2021, Keith et al Chem Rev 2021]

# Molecular Dynamics (MD) with ML



**Reducing necessary data by factor 1000+**

[Chmiela et al 2017, 2018, 2019, Noe et al. 2020, Sauceda et al 2021]

# Machine Learning meets DH

# (1) Machine learning and DH example

Alterations in historical manuscripts (Lassner et al. DHQ 2021)

# (2) Machine learning and DH example

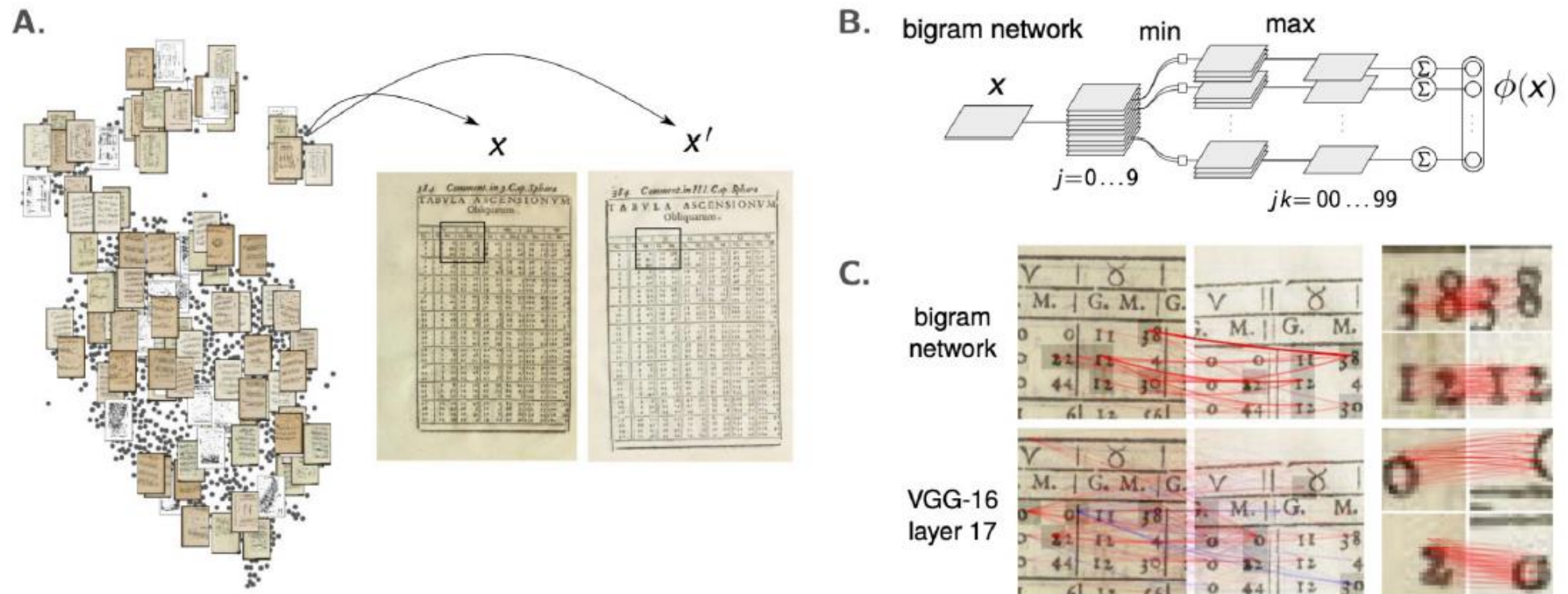Assessment of similarity between astronomical tables (Eberle et al. TPAMI 2020)



Fig. 11. **A.** Collection of tables from the Sphaera Corpus [29] from which we extract two tables with identical content. **B.** Proposed 'bigram network' supporting the table similarity model. **C.** BiLRP explanations of predicted similarities between the two input tables.

Will be explained in detail on Wednesday!

# (1) Methodologies

- In the sciences, scholars are used to statistical methods
- In the humanities, there is still a broader community that question statistical methods to be used in their field

# (2) Methodologies that humanities scholars might criticise

"We train a model on books and their amazon star rating to predict how good an unseen book is"

- It is not well defined what "good" means and how it relates to amazon star rating
- What is a book, how does that relate to literature or literariness
- The model gives an end-to-end answer

# (3) Methodologies that humanities scholars might accept

1. Falsification:

There is a general consensus in the field about a certain topic and our experiment presents evidence against that. For example: A book that was published anonymously has always been considered to be written by one person and a stylometric analysis now suggests that it is more similar to works by another person

2. Using ML as an intermediate transformation step:

Instead of giving an end-to-end solution, ML could be used to transform the input data such that it is more accessible to the scholar. OCR, Named Entity Recognition, Object Detection in paintings to make collections searchable by objects, etc. Here, the analysis is performed by an expert after the transformation

(Ramsay 2011)

# ML with few data

- The success of ML in recent years is in part accounted to increasing data set sizes
- In DH projects, we usually don't have a lot of data (and it is also not always possible to acquire new data - you cannot ask Goethe to write another novel)

We should therefore:

- Try to use the available data as efficiently as possible
- Try to use transfer learning
- If there is a lot of unlabelled data, weak labels might help

# Pros and cons of transfer learning

Pro: larger data sets can be used for pretraining

Pro: an unsupervised objective can be used for pretraining

Con: the more different the pretraining data is from the project data the less it will be transferable

Con: Large pretraining data may contain unwanted biases (more on this on thursday!)

# Collaboration between ML and DH from the ML perspective

- Acknowledge humanists careful approach to sources/data
- Try to understand humanists definitions (Example: ask a literary scholar "What is the text?")
- Managing expectations of ML solutions

## Conclusion

- need for opening the blackbox …

- understanding nonlinear models is essential for Sciences, DH & AI

- ML: a tool for gaining **<span style="color:red">insight</span>**

- **DH holds formidable challenges for ML**

- **teaching a new bi-lingual generation is necessary**