

ABOUT DATA FROM A MACHINE LEARNING PERSPECTIVE

Oriol Pujol



UNIVERSITAT DE
BARCELONA

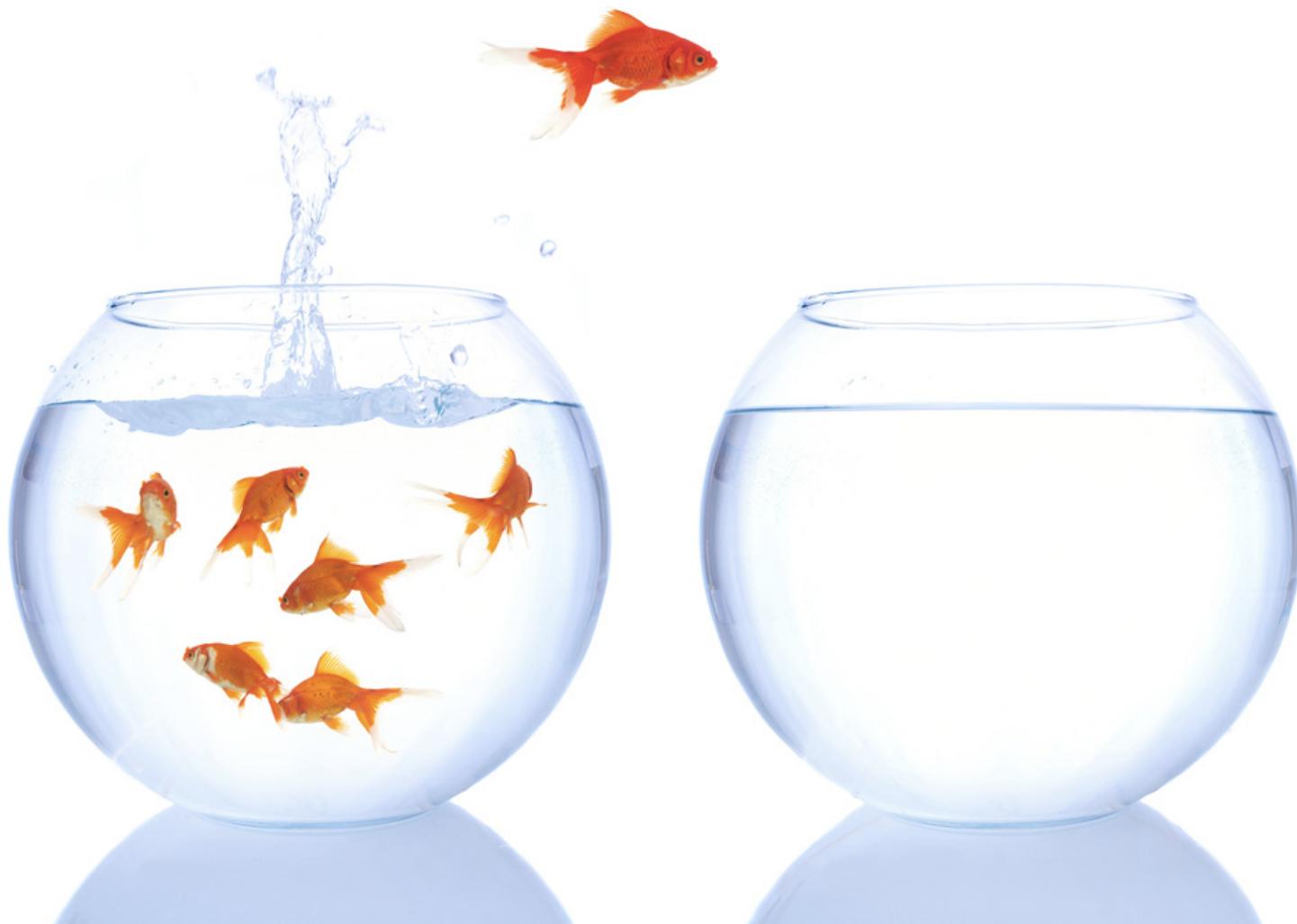


TWO FLAVOURS AND ONE PRESENTATION

Share experiences in openness from the machine learning community

To propose potential solutions to the problem of sharing sensible data

A COMMUNITY SHIFT



THE ROLE OF DATA IN MACHINE LEARNING

- Machine learning is sometimes called learning from data. Thus we are talking about building/coding algorithms that learn from data.
- Data is the key component in most research, but it is obviously crucial in machine learning.

EXAMPLES OF TASKS IN MACHINE LEARNING

Automatic captioning



"little girl is eating piece of cake."



"baseball player is throwing ball
in game."



"woman is holding bunch of
bananas."

CURRENT MACHINE LEARNING RESEARCH FOCUS

- Improving results in terms of generalization, robustness, speed, and interpretability.
- Comparison with other techniques is fundamental.
- Many different data sets are needed to assess each property we are tackling in the research, e.g. we need large datasets for speed and big data.
- Sometimes, their application largely depends on the domain, e.g., financial products forecasting, targeted gene identification.

AS A RESEARCHER

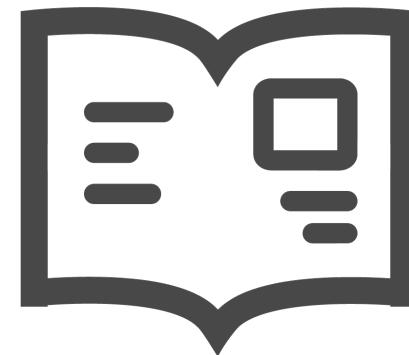
We pursue to contribute to advance science: reproducibility of experimentation that helps to eliminate individual biases.

DISEMINATION: REPORTING PROCESS, CLAIMING AUTHORSHIP

QUALITY: COMMUNITY VALIDITY ASSESSMENT, PEER REVIEW
(KPI USED BY ORGANISATIONS FOR RESEARCH QUALITY)

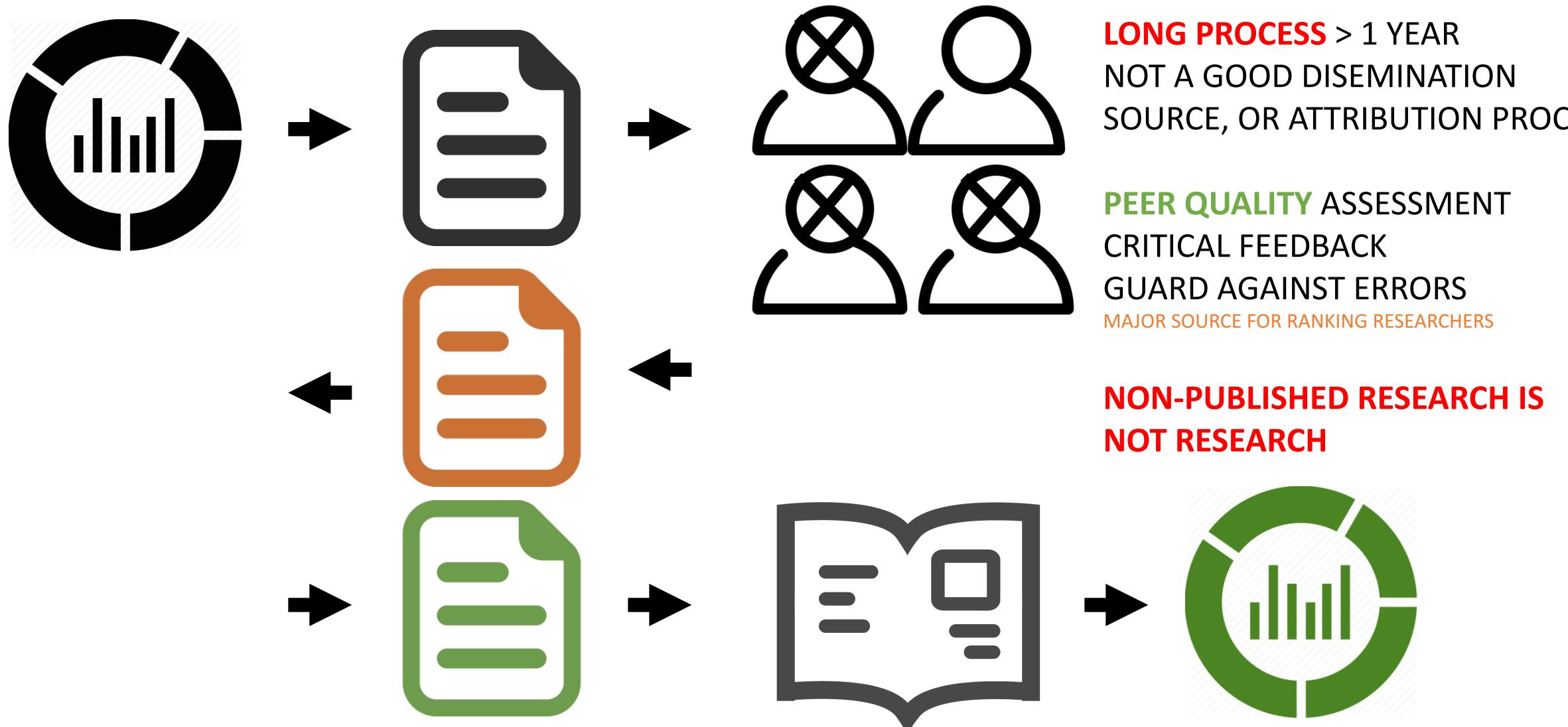


FOR THIS REASON WE USE PUBLISHING

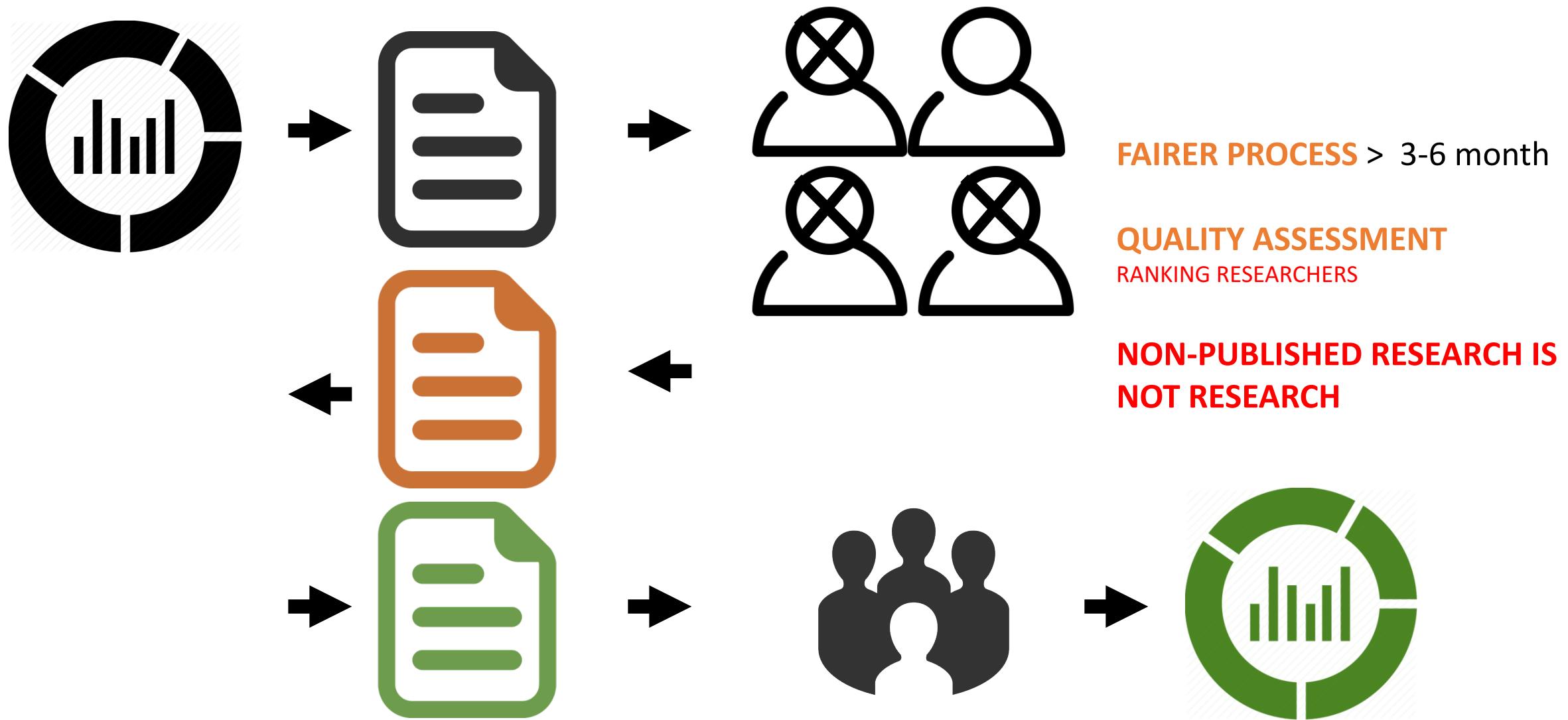


DISSEMINATION
QUALITY ASSESSMENT

FOR THIS REASON WE USE PUBLISHING



FOR THIS REASON WE USE PUBLISHING



THE NIPS CASE: WHEN CONFERENCE PAPERS ARE OLD NEWS



IN MACHINE LEARNING 3-6 MONTHS IS TOO LONG

DECEMBER 5TH – 10TH, 2016

NIPS Proceedings^β

Books

2016

Conditional Image Generation with PixelCNN Decoders

Part of: [Advances in Neural Information Processing Systems 29 \(NIPS 2016\) pre-proceedings](#)

[PDF] [BibTeX] [Supplemental] [Reviews]

Authors

- [Aaron van den Oord](#)
- [Nal Kalchbrenner](#)
- [Lasse Espeholt](#)
- [koray kavukcuoglu](#)
- [Oriol Vinyals](#)
- [Alex Graves](#)

THE NIPS CASE: WHEN CONFERENCE PAPERS ARE OLD NEWS

CORNELL UNIVERSITY LIBRARY
Cornell University Library

We gratefully acknowledge support from the Simons Foundation and Stockholm University

arXiv.org > cs > arXiv:1606.05328

Search or Article ID inside arXiv All papers Broaden your search using Semantic Scholar

(Help | Advanced search)

Computer Science > Computer Vision and Pattern Recognition

Conditional Image Generation with PixelCNN Decoders

Aaron van den Oord, Nal Kalchbrenner, Oriol Vinyals, Lasse Espeholt, Alex Graves, Koray Kavukcuoglu

(Submitted on 16 Jun 2016 (v1), last revised 18 Jun 2016 (this version, v2))

This work explores conditional image generation with a new image density model based on the PixelCNN architecture. The model can be conditioned on any vector, including descriptive labels or tags, or latent embeddings created by other networks. When conditioned on class labels from the ImageNet database, the model is able to generate diverse, realistic scenes representing distinct animals, objects, landscapes and structures. When conditioned on an embedding produced by a convolutional network given a single image of an unseen face, it generates a variety of new portraits of the same person with different facial expressions, poses and lighting conditions. We also show that conditional PixelCNN can serve as a powerful decoder in an image autoencoder. Additionally, the gated convolutional layers in the proposed model improve the log-likelihood of PixelCNN to match the state-of-the-art performance of PixelRNN on ImageNet, with greatly reduced computational cost.

Subjects: Computer Vision and Pattern Recognition (cs.CV); Learning (cs.LG)
Cite as: arXiv:1606.05328 [cs.CV]
(or arXiv:1606.05328v2 [cs.CV] for this version)

Submission history

From: Aäron van den Oord [view email]
[v1] Thu, 16 Jun 2016 19:40:56 GMT (3016kb,D)
[v2] Sat, 18 Jun 2016 15:44:24 GMT (3016kb,D)

Download:

- PDF
- Other formats

(license)

Current browse context:
cs.CV
< prev | next >
new | recent | 1606

Change to browse by:
cs
 cs.LG

References & Citations
• NASA ADS

DBLP – CS Bibliography
listing | bibtex
Aäron van den Oord
Nal Kalchbrenner
Oriol Vinyals
Lasse Espeholt
Alex Graves
...

24/7 PUBLISHING CYCLE



IMMEDIACY, 24/7 PUBLISHING CYCLE

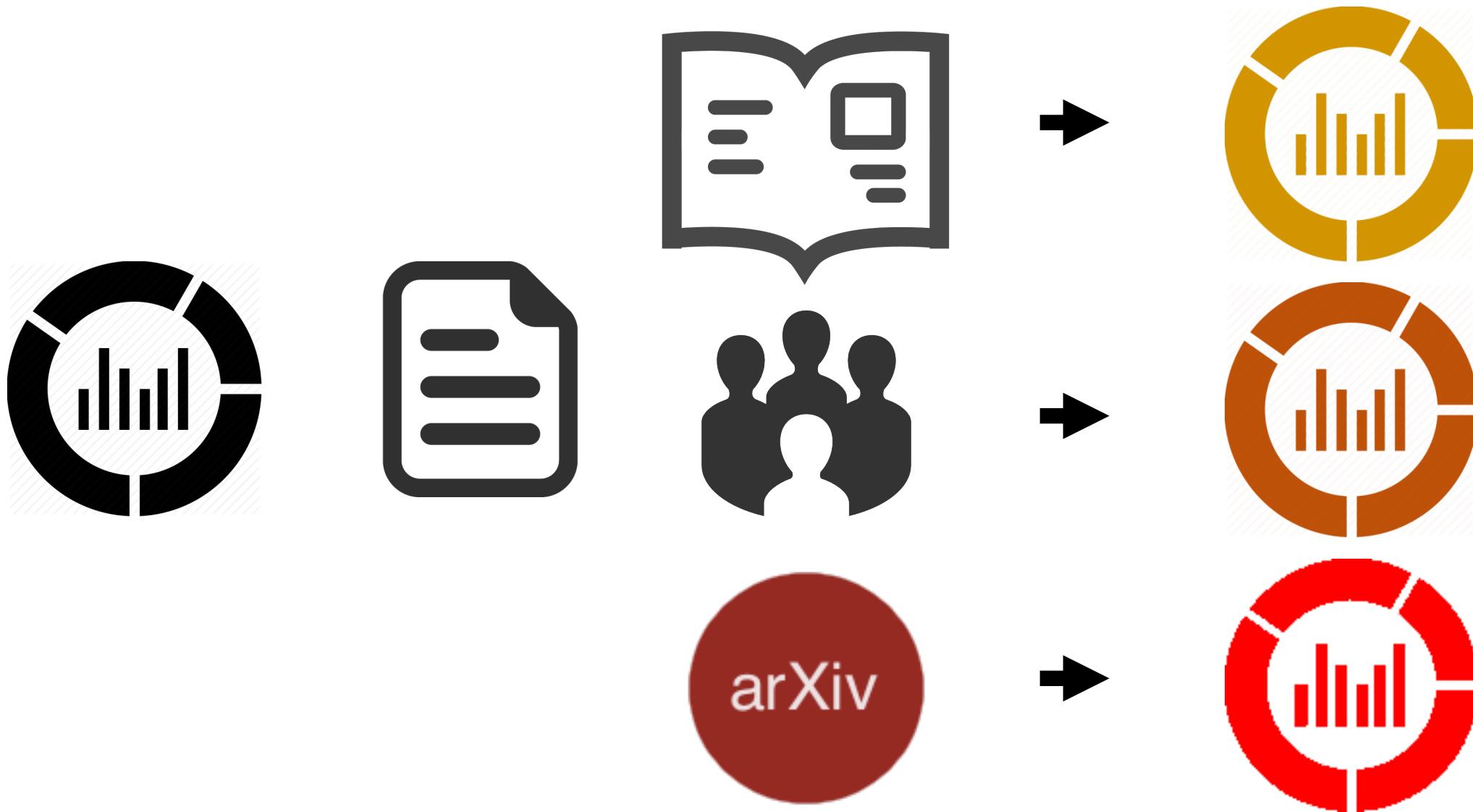
IMMEDIATE ATTRIBUTION

NO QUALITY ASSESSMENT

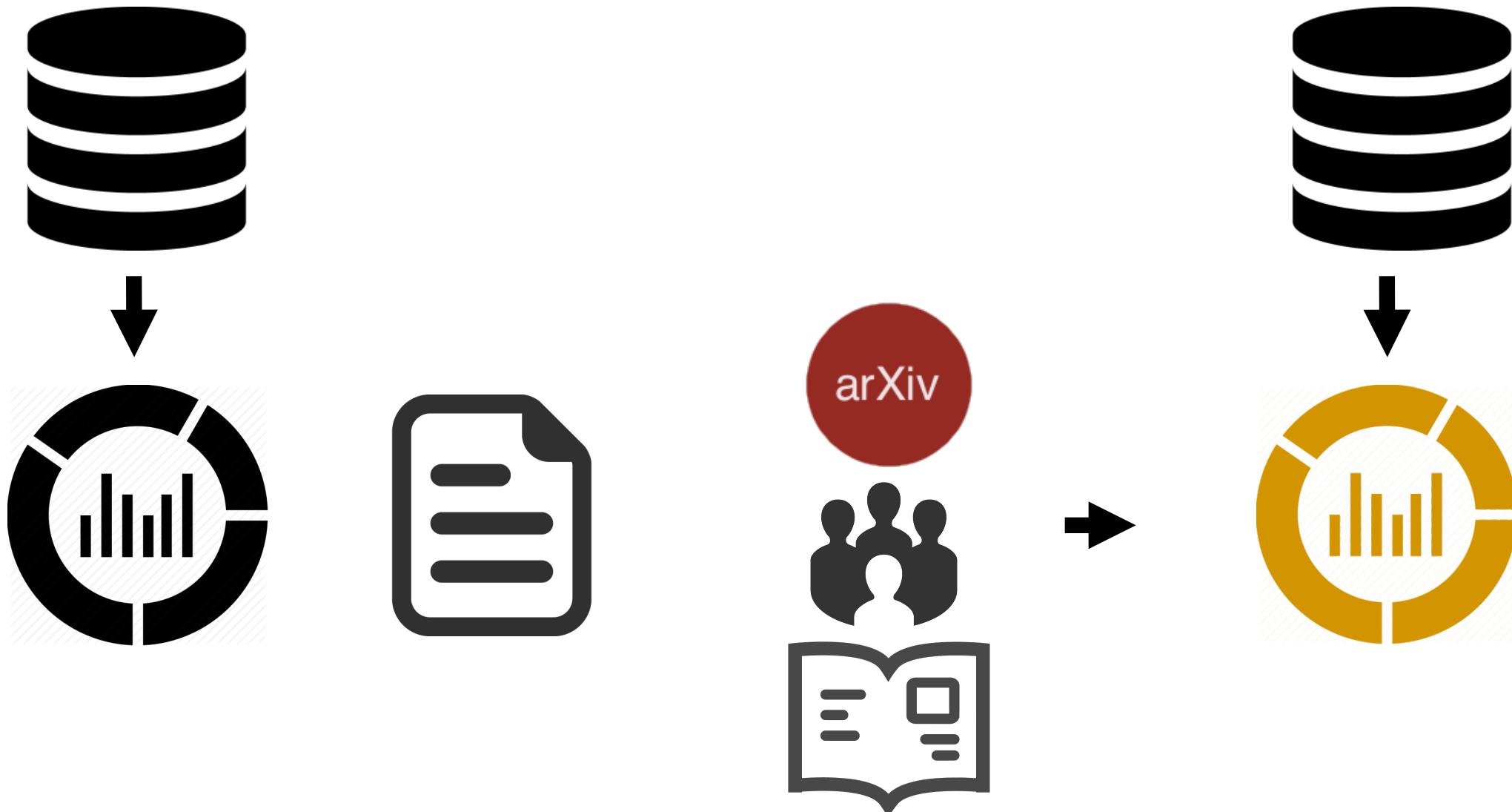
RANKING RESEARCHERS

HOW TO DISTINGUISH GOOD RESEARCH FROM FRAUD

HOWEVER, IN REALITY



TOWARDS REPRODUCIBILITY: WE NEED DATA



AN EXAMPLE: PLOS ONE

FOCUS SHIFT:

METHODOLOGICAL
CORRECTNESS

DATA

The screenshot shows the PLOS ONE homepage with a "TENTH ANNIVERSARY" banner. Below the banner, there are three tabs: "In the News" (highlighted), "Recent", and "Most Viewed". The main content area features a research article titled "Reconstruction of the Cortical Maps of the Tasmanian Tiger and Comparison to the..." by Gregory S. Berns, Ken W. S. Ashwell. It includes two brain scan images labeled "THYLACINE 1" and "THYLACINE 2". Below the article are two graphs labeled "a." and "b." Both graphs plot "Sensitivity" (y-axis, 0.0 to 1.0) against "1 - Specificity" (x-axis, 0.0 to 1.0). Graph "a." compares "NSE at 48 hours" (solid line, AUC = 0.83) and "NSE at 72 hours" (dashed line, AUC = 0.85). Graph "b." compares "NSE at 48 hours" (solid line, AUC = 0.86) and "Geometrical area under the NSE-curve at 24 to 72 hours" (dashed line, AUC = 0.86).

PLOS | ONE
TENTH ANNIVERSARY

In the News Recent Most Viewed

THYLACINE 1 THYLACINE 2

Reconstruction of the Cortical Maps of the Tasmanian Tiger and Comparison to the...

Gregory S. Berns, Ken W. S. Ashwell

a.

b.

Sensitivity

1 - Specificity

DLG MGN LP VPLM VAVL

n = 397

— NSE at 48 hours AUC = 0.83
— NSE at 72 hours AUC = 0.85

n = 370

— NSE at 48 hours AUC = 0.86
— Geometrical area under the NSE-curve at 24 to 72 hours AUC = 0.86

Single versus Serial Measurements of Neuron-Specific Enolase and Prediction of Poor...

Sebastian Wiberg, Christian Hassager, [...], Jesper Kjaergaard

AUC = 0.86

MACHINE LEARNING DATASET REPOSITORIES


Machine Learning Repository
Center for Machine Learning and Intelligent Systems

About [Citation Policy](#) [Donate a Data Set](#) [Contact](#)
Loading
[View ALL Data Sets](#)

Welcome to the UC Irvine Machine Learning Repository!

We currently maintain 360 data sets as a service to the machine learning community. You may [view all data sets](#) through our searchable interface. Our [old web site](#) is still available, for those who prefer the old format. For a general overview of the Repository, please visit our [About page](#). For information about citing data sets in publications, please read our [citation policy](#). If you wish to donate a data set, please consult our [donation policy](#). For any other questions, feel free to [contact the Repository librarians](#). We have also set up a [mirror site](#) for the Repository.

Supported By:  In Collaboration With:  [Rexa.info](#)
• Research • People • Connections

Latest News: <p>2013-04-04: Welcome to the new Repository admins Kevin Bache and Moshe Lichman! 2010-03-01: Note from donor regarding Netflix data 2009-10-16: Two new data sets have been added. 2009-09-14: Several data sets have been added. 2008-07-23: Repository mirror has been set up. 2008-03-24: New data sets have been added! 2007-06-25: Two new data sets have been added: UJI Pen Characters, MAGIC Gamma Telescope</p>	Newest Data Sets: <p>2016-11-23:  NIPS Conference Papers 1987-2015 2016-11-16:  Amazon book reviews 2016-08-14:  Dota2 Games Results 2016-08-05:  Facebook metrics</p>	Most Popular Data Sets (hits since 2007): <p>1189468:  Iris 816286:  Adult 621858:  Wine 534408:  Car Evaluation</p>
---	--	--

MACHINE LEARNING DATASET REPOSITORIES

Facebook metrics Data Set

Download: [Data Folder](#) [Data Set Description](#)

Abstract: Facebook performance metrics of a renowned cosmetic's brand Facebook page.

Data Set Characteristics:	Multivariate	Number of Instances:	500	Area:	Business
Attribute Characteristics:	Integer	Number of Attributes:	19	Date Donated	2016-08-05
Associated Tasks:	Regression	Missing Values?	N/A	Number of Web Hits:	7110

Source:

Created by: Sárgio Moro, Paulo Rita and Bernardo Vala (ISCTE-IUL) @ 2016

Data Set Information:

The data is related to posts' published during the year of 2014 on the Facebook's page of a renowned cosmetics brand. This dataset contains 500 of the 790 rows and part of the features analyzed by Moro et al. (2016). The remaining were omitted due to confidentiality issues.

Attribute Information:

It includes 7 features known prior to post publication and 12 features for evaluating post impact (see Tables 2 and 3 from Moro et al., 2016 - complete reference in the 'Citation Request')

Relevant Papers:

(Moro et al., 2016) Moro, S., Rita, P., & Vala, B. (2016). Predicting social media performance metrics and evaluation of the impact on brand building: A data mining approach. *Journal of Business Research*, 69(9), 3341-3351.

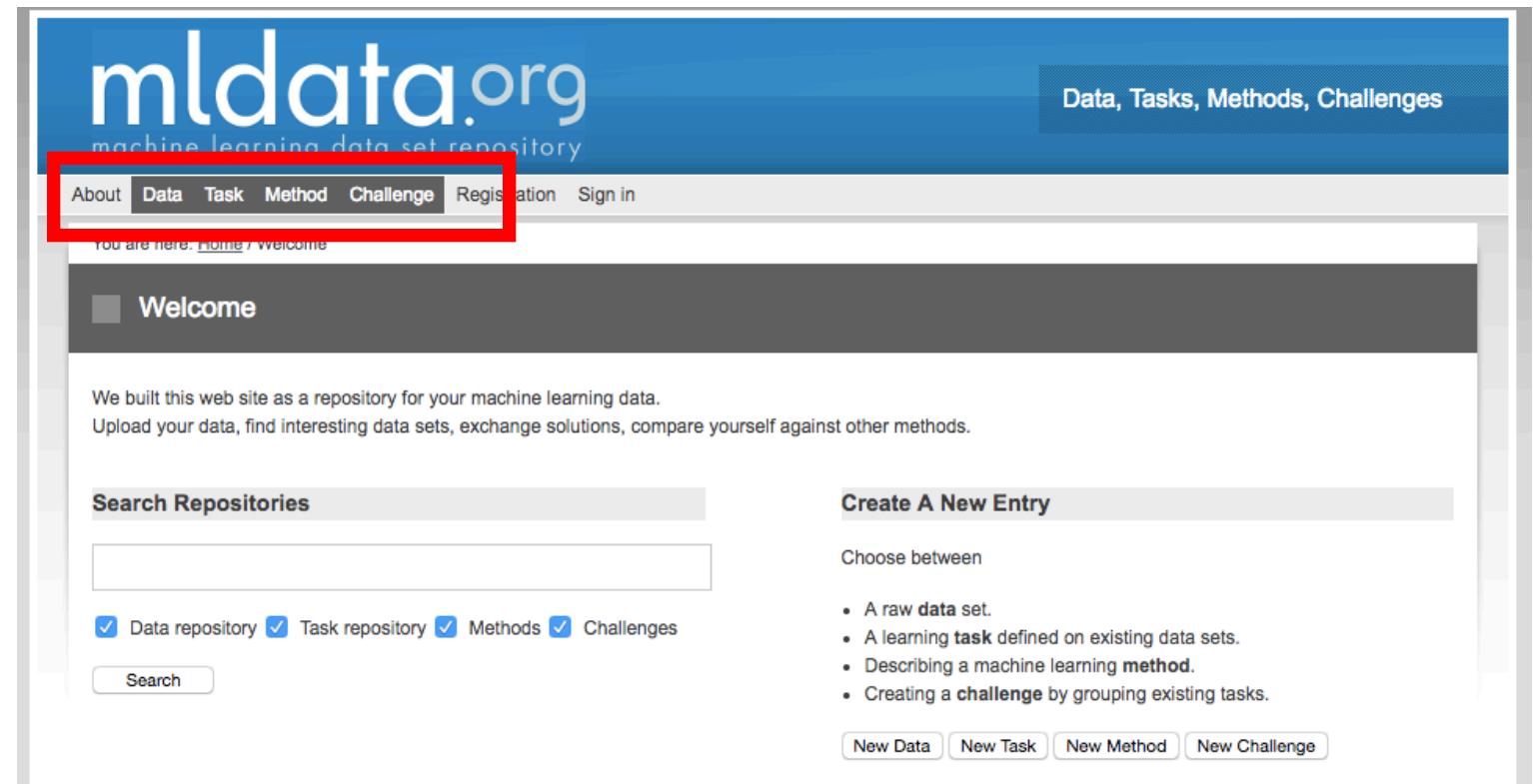
MACHINE LEARNING DATASET REPOSITORIES

Data Sets - Raw data as a collection of similarly structured objects.

Material and Methods - Descriptions of the computational pipeline.

Learning Tasks - Learning tasks defined on raw data.

Challenges - Collections of tasks which have a particular theme.



WHEN DESCRIPTION AND DATA IS NOT ENOUGH



U.S. SECURITIES AND EXCHANGE COMMISSION

- *The U.S. Securities and Exchange Commission (SEC) 2010 proposal, financial products must include Python code.*

SECURITIES AND EXCHANGE COMMISSION

17 CFR Parts 200, 229, 230, 232, 239, 240, 243 and 249

Release Nos. 33-9117; 34-61858; File No. S7-08-10

RIN 3235-AK37

ASSET-BACKED SECURITIES

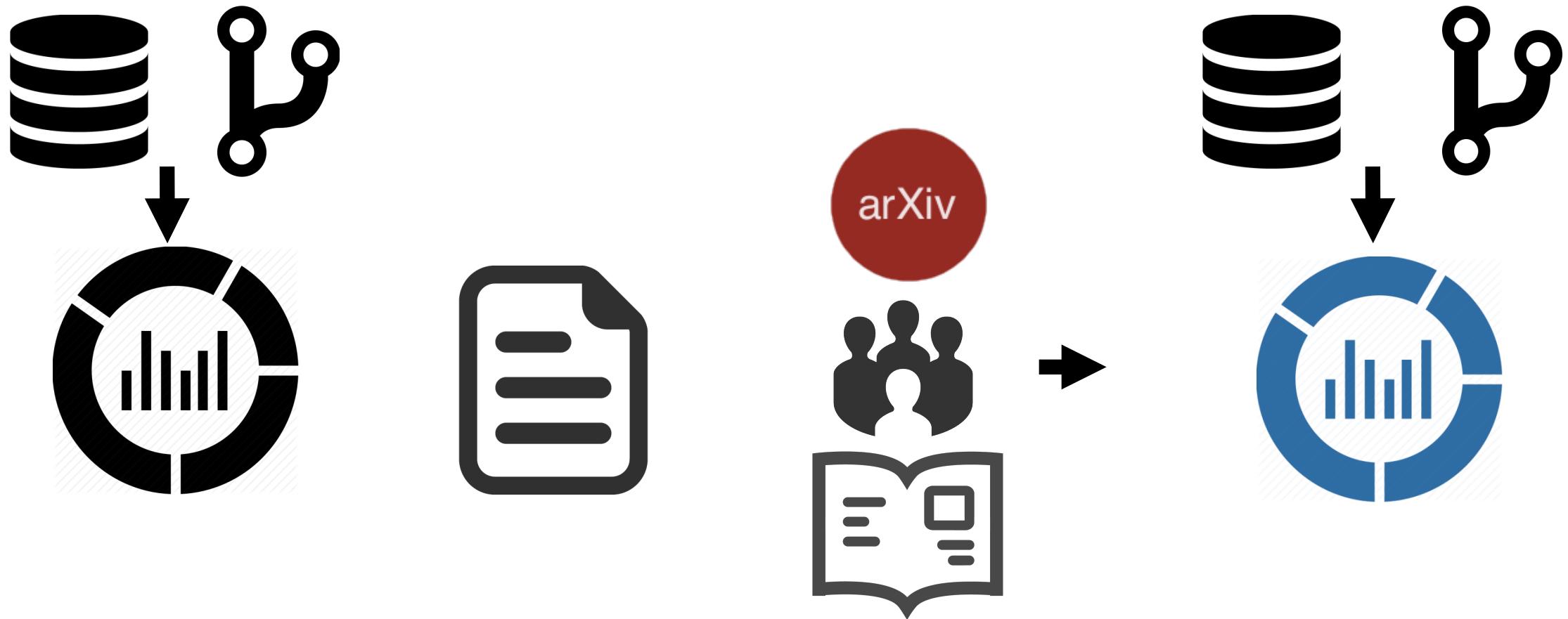
AGENCY: Securities and Exchange Commission

ACTION: Proposed rule.

SUMMARY: We are proposing significant revisions to Regulation AB and other rules regarding the offering process, disclosure and reporting for asset-backed securities. Our proposals would revise filing deadlines for ABS offerings to provide investors with more time to consider transaction-specific information, including information about the pool assets. Our proposals also would repeal the current credit ratings references in shelf eligibility criteria for asset-backed issuers and establish new shelf eligibility criteria that would include, among other things, a requirement that the sponsor retain a portion of each tranche of the securities that are sold and a requirement that the issuer undertake to file Exchange Act reports on an ongoing basis so long as its public securities are outstanding. We also are proposing to require that, with some exceptions, prospectuses for public offerings of asset-backed securities and ongoing Exchange Act reports contain specified asset-level information about each of the assets in the pool. The asset-level information would be provided according to proposed standards and in a tagged data

format using eXtensible Markup Language (XML). In addition, we are proposing to require, along with the prospectus filing, the filing of a computer program of the contractual cash flow provisions expressed as downloadable source code in Python, a

WHEN DESCRIPTION AND DATA IS NOT ENOUGH



GITHUB PROJECTS

The screenshot shows a GitHub repository page. At the top, there's a navigation bar with links for Personal, Open source, Business, Explore, Pricing, Blog, Support, and a search bar. Below the navigation is the repository name 'oriolpujol / pyday_BCN2016_learning_representations'. To the right of the name are buttons for Watch (3), Star (2), and Fork (0). Below the repository name, there are tabs for Code (selected), Issues (0), Pull requests (0), Projects (0), Pulse, and Graphs.

This is the notebook used for pyDay BCN 2016: XXL Meetup about "Learning Representations using Deep Learning with TensorFlow"

Key statistics shown: 3 commits, 1 branch, 0 releases, 1 contributor, and BSD-3-Clause license.

Branch dropdown: master ▾ | New pull request | Find file | Clone or download ▾

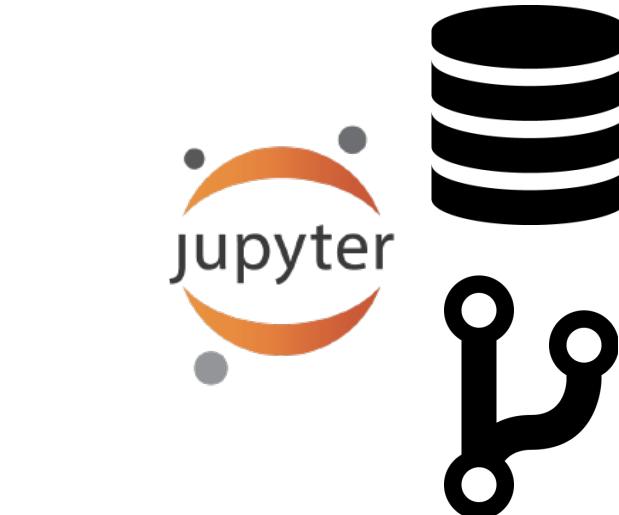
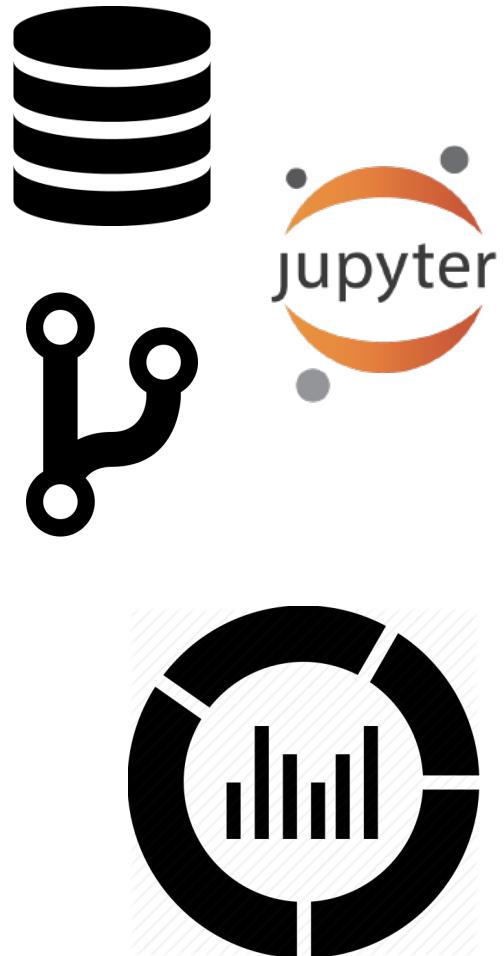
Commit history:

File	Message	Date
LICENSE	Update and rename LICENSE.md to LICENSE	4 months ago
NN.jpg	Add files via upload	4 months ago
bandwagon.jpg	Add files via upload	4 months ago
neuron.jpg	Add files via upload	4 months ago
pyBCN_learning_representations_deep_learn...	Bug fixes and tqdm progress bar added	4 months ago

BUT ... DEVIL IS IN THE DETAILS



DOCUMENTING THE DETAILS



JUPYTER NOTEBOOKS

EXECUTABLE CODE
LATEX
MARKDOWN
HTML+CSS
HYPERLINKS
EMBEDDED FRAMES
... and more.

jupyter pill8_Stochastic_Subgradient_Methods_master Last Checkpoint: 10/21/2016 (autosaved) Python 2 O

File Edit View Insert Cell Kernel Help

Cell Toolbar

```
16 print z.shape
17 z.shape=sz
18 plt.imshow(z, interpolation='bilinear', origin='lower', extent=(-2,2,-2,2),alpha=0.3, vmin=-10, vmax=10)
19 plt.contour(XX,YY,z,[0,1,2,3,4,5])
20 fig = plt.gcf()
21 fig.set_size_inches(9,9)
22 plt.plot(np.array(wpath)[0,0],np.array(wpath)[0,1],'ro')
23 plt.plot(np.array(wpath)[:,0],np.array(wpath)[:,1])
24 plt.plot(np.array(wpath)[-1,0],np.array(wpath)[-1,1],'co')
```

Adam

Adam combines both momentum and adagrad.

$$m^{(k+1)} = \gamma_1 m^{(k)} - (1 - \gamma_1)g^{(k)}$$
$$D^{(k+1)} = \gamma_2 D^{(k)} + (1 - \gamma_2)\text{diag}(g^{(k)})^2$$
$$x^{(k+1)} = x^{(k)} + \alpha_k (D^{(k+1)} + \epsilon I)^{-1/2} m^{(k+1)}$$

Usual values of γ_1 are 0.9 and γ_2 are 0.999, ϵ can take very small values.

EXERCISE Code the Adam according to the former equations.

PUTTING ALL TOGETHER

Competitions About Categories ▾ Search

GitXiv

Register Sign In

Post

Receive the best of GitXiv right in your inbox.

Your Email

Get Newsletter



Collaborative Open Computer Science

View : Top New Best Daily



Deep Probabilistic Programming

Edward, a Turing-complete probabilistic programming language

BAYESIAN DEEP LEARNING (DL)

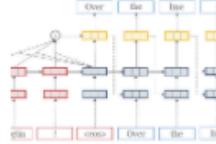
1



Christos Iraklis Tsatsoulis 1 point 5 hours ago 0 Comments



0



OpenNMT: Open-Source Toolkit for Neural Machine Translation

Industrial-strength, open-source neural machine translation system

NATURAL LANGUAGE PROCESSING (NLP) RECURRENT NEURAL NETWORKS (RNN)

2



Christos Iraklis Tsatsoulis 1 point 6 days ago 0 Comments



0



Improved Texture Networks: Quality & Diversity in Feed-forward Stylization & Texture Synthesis

Introducing Instance Normalization for better quality, entropy loss for diversity.

COMPUTER VISION CONVOLUTIONAL NEURAL NETWORKS (CNN) GENERATIVE

3



Dmitry Ulyanov 2 points 6 days ago 0 Comments



0



SalGAN: Visual Saliency Prediction with Generative Adversarial Networks

Adversarial training improves a binary cross-entropy loss in most saliency metrics

ADVERSARIAL NETWORKS COMPUTER VISION CONVOLUTIONAL NEURAL NETWORKS (CNN) DEEP LEARNING (DL)

Xavier Giró-i-Nieto 2 points 13 days ago 0 Comments



This paper introduces the usage of generative adversarial networks (GANs) for visual saliency prediction. In this context, training is driven by two agents. First, the Generator that creates a synthetic sample matching the data distribution modelled by a training data set; second, the Discriminator, that distinguishes between a real sample drawn directly from the training data set and one created by the generator.

arXiv

We introduce SalGAN, a deep convolutional neural network for visual saliency prediction trained with adversarial examples. The first stage of the network consists of a generator model whose weights are learned by back-propagation computed from a binary cross entropy (BCE) loss over downsampled versions of the saliency maps. The resulting prediction is processed by a

GitHub

Our paper presents two convolutional neural networks, one corresponds to the Generator (Saliency Prediction Network) and the other is the Discriminator for the adversarial training. To compute saliency maps only the Generator is needed. SalGAN is implemented in Lasagne, which at its time is developed over Theano.

Links

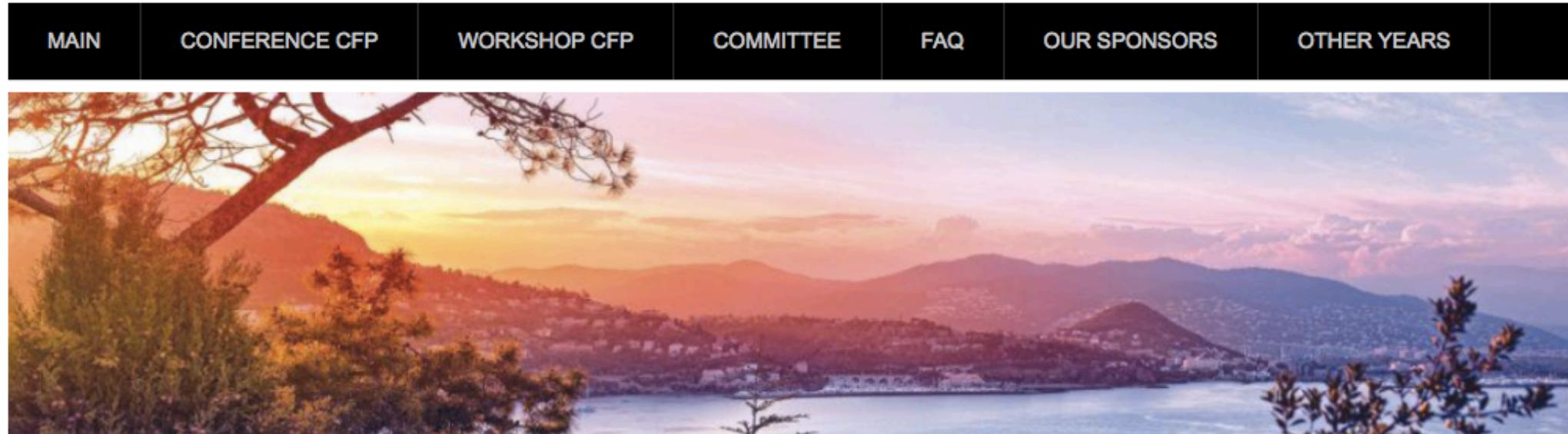
- <https://imatge-upc.github.io/saliency-salgan-2017/>.
- <http://www.slideshare.net/xavigiro/salgan-visual-saliency-prediction-with-generative-adversarial-networks>.

BOTTOM LINE

1. SCIENCE IS MORE THAN A PUBLICATION
2. AT ITS CORE LIES REPRODUCIBILITY...THIS INVOLVES DATA
3. BUT DATA IS NOT EVEN ENOUGH... CODE AND PRECISE EXPERIMENT
REPRODUCIBILITY
4. DATA AND PROGRAMMING LITERACY IS NEEDED

BUT ... SOMETHING ELSE ABOUT QUALITY ASSESSMENT: OPEN REVIEW CASE

ICLR 2017



5th International Conference on Learning Representations

Overview

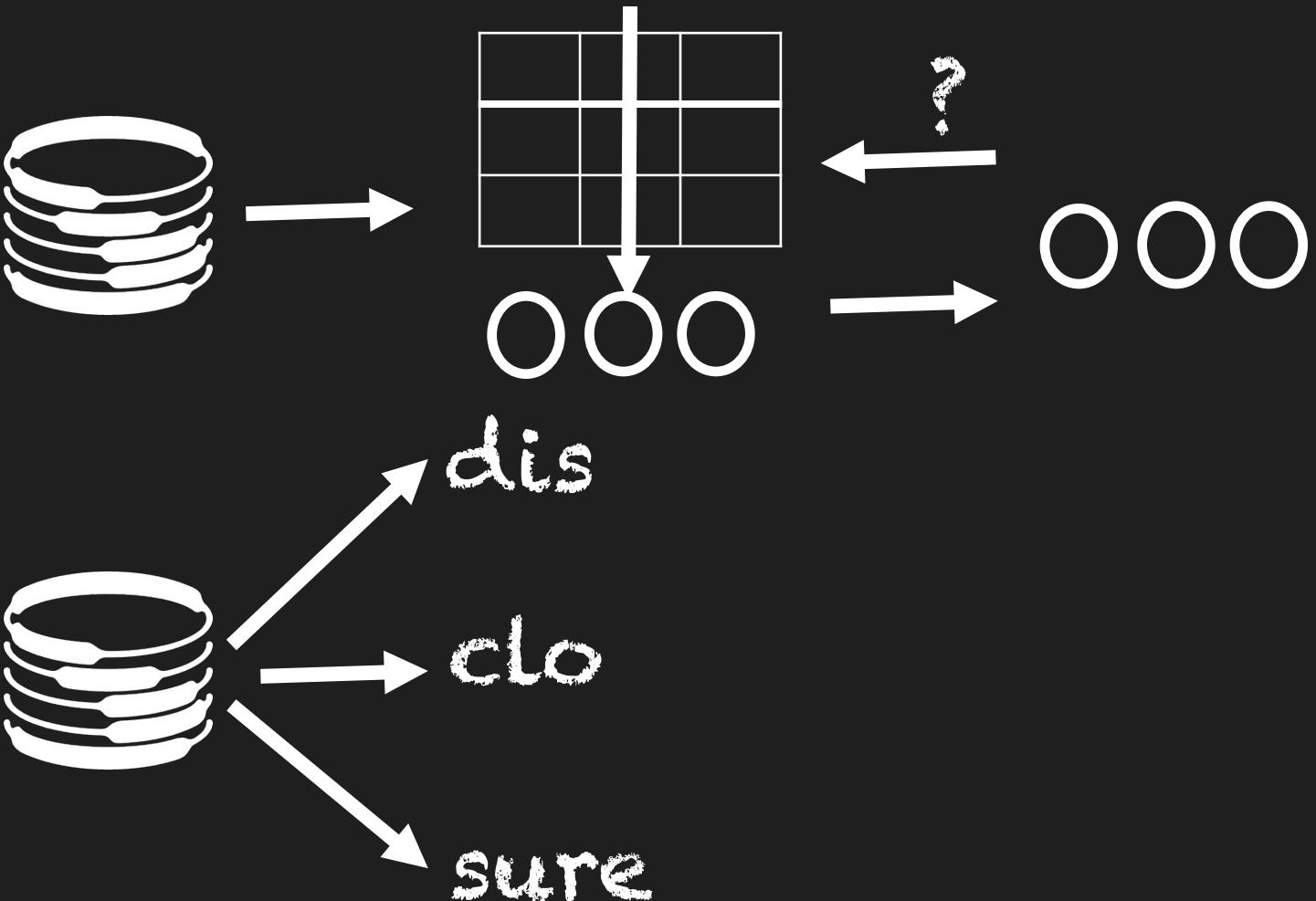
THE RULES

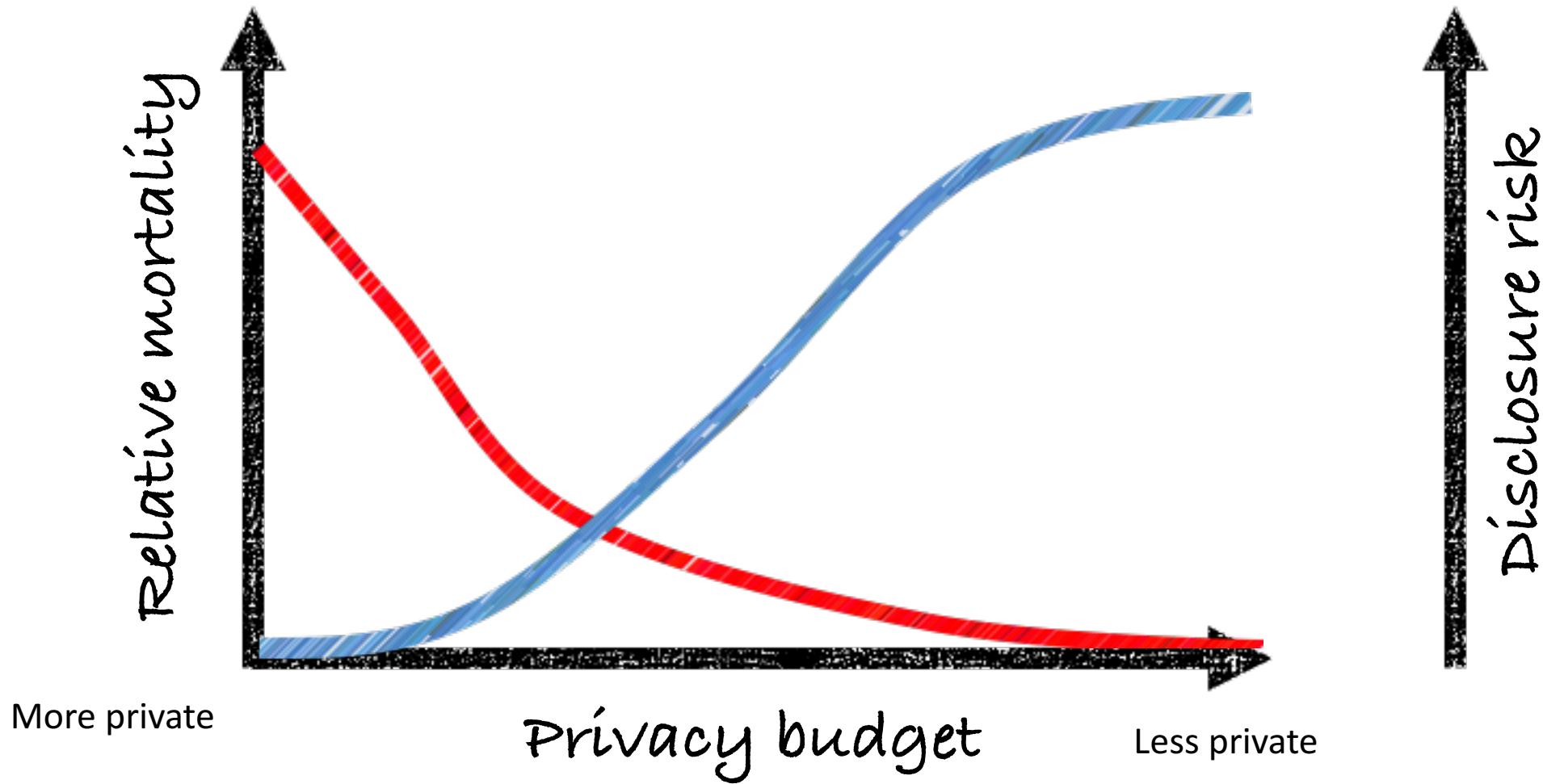
- Submissions posted on arXiv before being submitted to the conference
- Program committee assigns blind reviewers (as always) marked as designated reviewers
- Anyone can ask the program chairs for permission to become an anonymous designated reviewer (open bidding).
- Open commenters will have to use their real names, linked with their Google Scholar profiles.
- Papers that are presented in the workshop track or are not accepted will be considered non-archival, and may be submitted elsewhere (modified or not), although the ICLR site will maintain the reviews, the comments, and the links to the arXiv versions

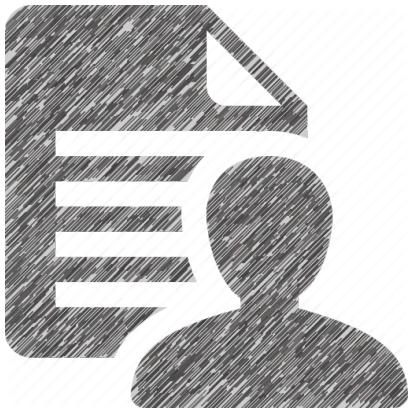
SOME POINTS ABOUT OPEN REVIEW

- Reading the answers from others is valuable
- Rejected papers have influence
- May become a popularity contest

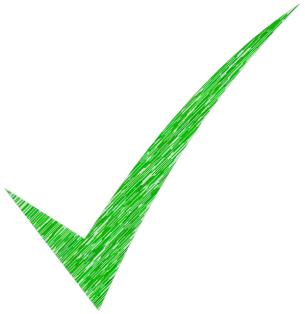
SHARING DATA IN THE AGE OF DEEP LEARNING

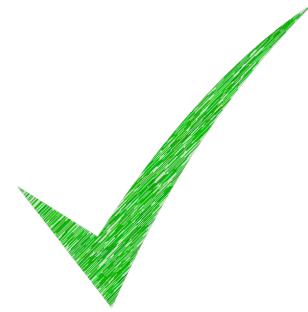
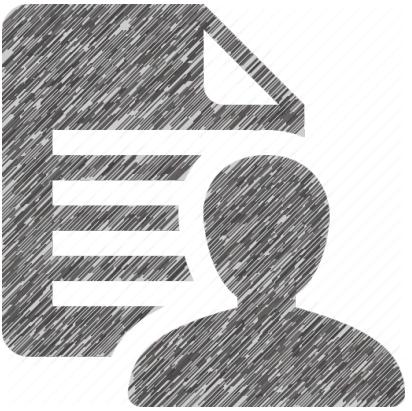


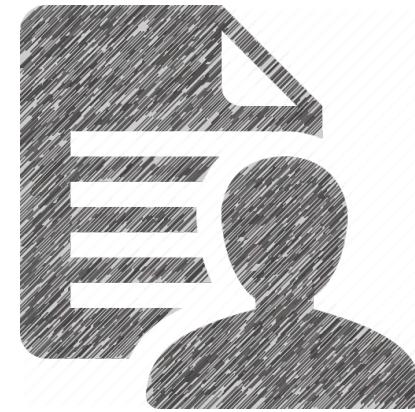
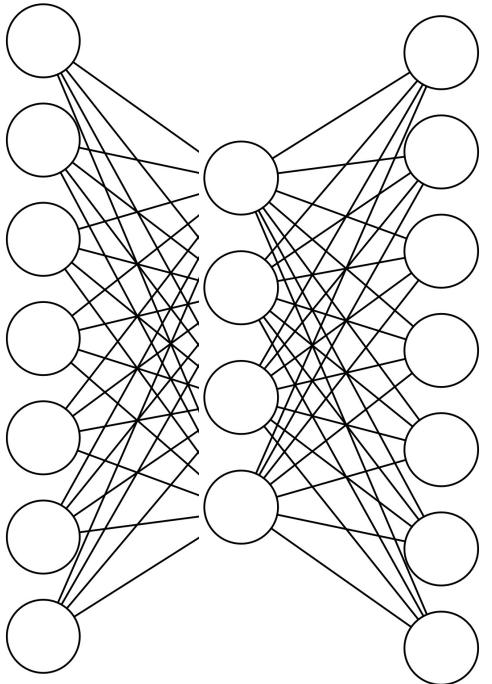
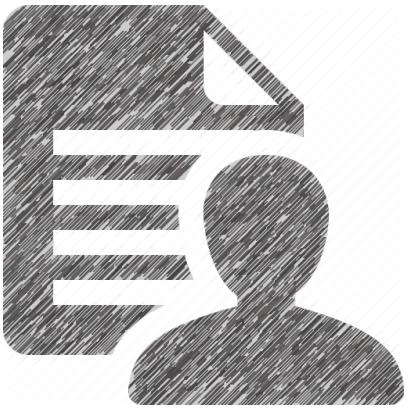


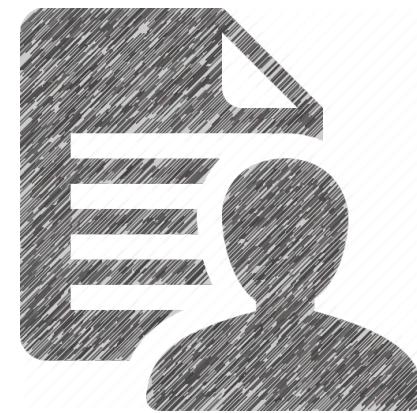
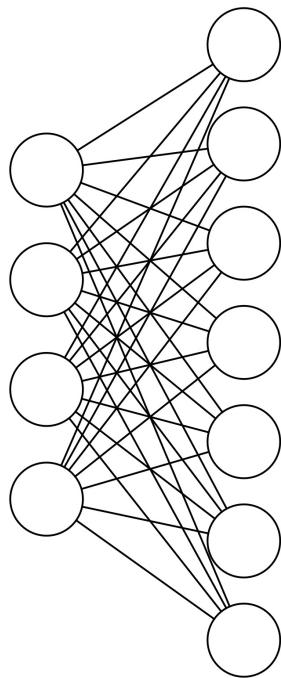
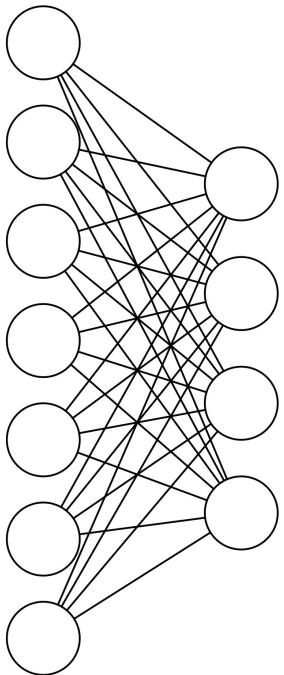
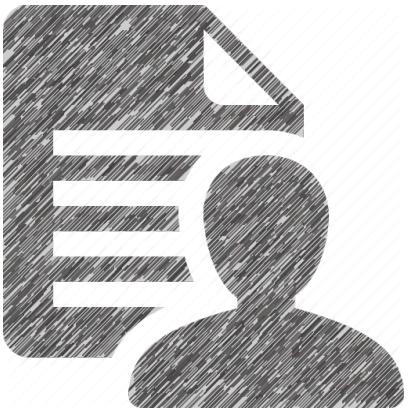


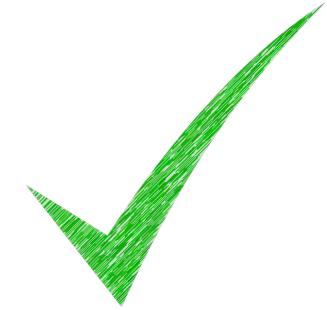
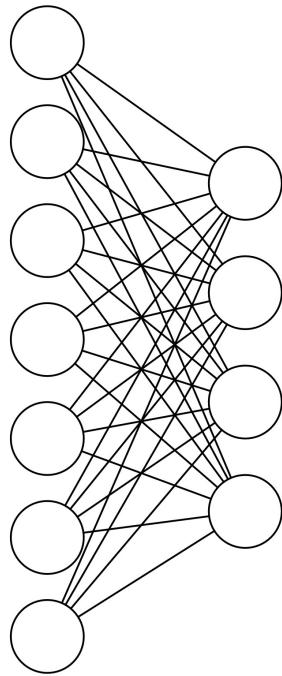
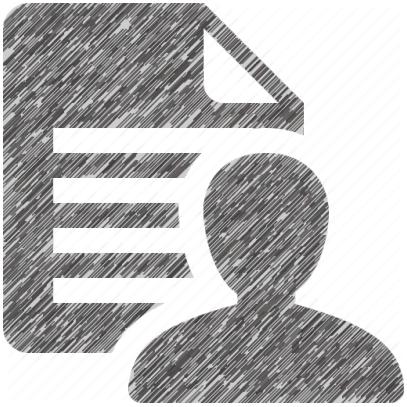
CLASSIFIER

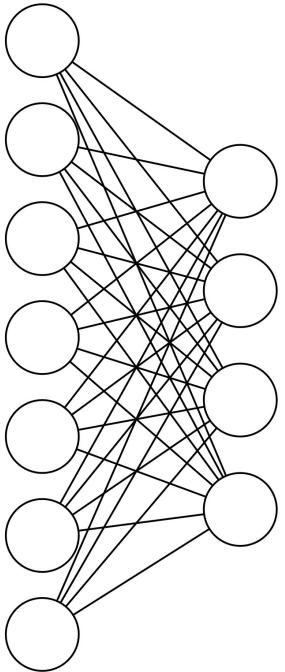
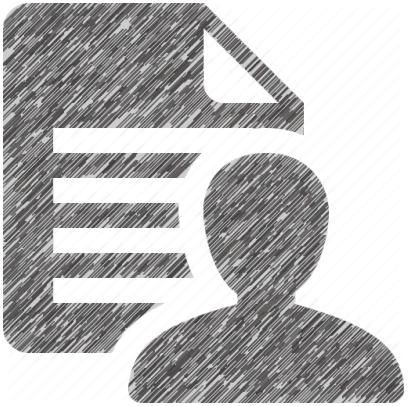




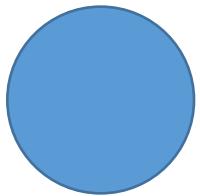




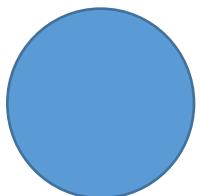




CONCLUSIONS



Research is based on reproducibility: documentation, data, code, and details.



Machine learning methods can help in the challenge for democratizing data.

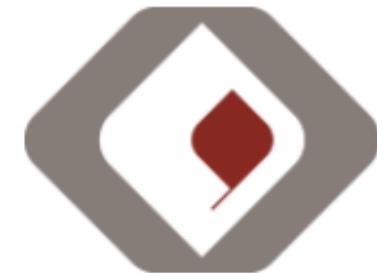
Thank you

oriol_pujol@ub.edu

@oriolpujolvila



UNIVERSITAT DE
BARCELONA



learn