



**OPEN SOURCE**  
LEADERSHIP SUMMIT

# **Introduction to the Community Data License Agreement**

**Nick Acosta, IBM**

# Introduction



# Agenda

<b>Part One - Introduction</b>	<b>01</b>	<b>Part Four - Example, pt. 1</b>	<b>26</b>
Agenda	03	Deep Learning Developers Dataset	27
Open Source Data	04	Deep Learning Frameworks	28
<b>Part Two - CDLA</b>	<b>10</b>	<b>Part Five - Example, pt. 2</b>	<b>33</b>
Overview	11	Double Pendulum Chaotic Dataset	34
Purpose	13	Learning beyond simulated physics	36
Computational Use	15		
 		<b>Part Six - Conclusion</b>	<b>37</b>
<b>Part Three - CDLA Versions</b>	<b>17</b>		
Versions	18		
Example	21		

# Why Open Source Data

[Data sharing uncovers five new...](#)

[https://www.nih.gov/news-events/news-releases/data-sharing-uncovers-five-new-risk-genes-alzheimers-disease](#)

National Institutes of Health

Turning Discovery Into Health

Search NIH

[NIH Employee Intranet](#)
[Staff Directory](#)
[En Espaol](#)

Health Information

Grants & Funding

News & Events

Research & Training

Institutes at NIH

About NIH

Home » News & Events » News Releases

NEWS RELEASES

Thursday, February 28, 2019

Data sharing uncovers five new risk genes for Alzheimer's disease

NIH-funded project includes largest sample to date for Alzheimer's gene association.

Analysis of genetic data from more than 94,000 individuals has revealed five new risk genes for Alzheimer's disease, and confirmed 20 known others. An international team of researchers also reports for the first time that mutations in genes specific to tau, a hallmark protein of Alzheimer's disease, may play an earlier role in the development of the disease than originally thought. These new findings support developing evidence that groups of genes associated with specific biological processes, such as cell trafficking, lipid transport, inflammation and the immune response, are "genetic hubs" that are an important part of the disease process. The study, which was funded in part by the National Institute on Aging (NIA) and other components of the National Institutes of Health, follows [results from 2013](#). It will be published online February 28, 2019 in the journal *Nature Genetics*.

Newly identified (red) and known (blue) genes linked to Alzheimer's disease spike in this table plotting results from genome-wide association analysis of 94,437 individuals with late onset Alzheimer's. *Kunkle et al and Nature Genetics.*

Institute/Center

National Institute on Aging (NIA)

Contact

Joe Balintfy

301-496-1752

Connect with Us

Subscribe to news releases

RSS Feed

"This continuing collaborative research into the genetic underpinnings of Alzheimer's is allowing us to dig deeper into the complexities of this devastating disease," said [Richard J. Hodes, M.D.](#), director of the NIA. "The size of this study provides additional clarity on the genes to prioritize as we continue to better understand and target ways to treat and prevent Alzheimer's."

The researchers, members of the International Genomic Alzheimer's Project (IGAP), analyzed both rare and common gene variants in 94,437 individuals with late onset Alzheimer's disease, the most common form of dementia in older adults. IGAP is made up of four centers in the United States and Europe that have been collaborating since 2011 on genome-wide association studies (GWAS).

# The National Institute on Aging Genetics of Alzheimer's Disease Data Storage

## NIAGADS DATA REQUEST DOCUMENTS

DOCUMENT NAME	UPLOAD DATE
<a href="#">Data Distribution Agreement</a>	2018-10-29
<a href="#">NIAGADS Application Instructions</a>	2016-09-09
<a href="#">NIAGADS Renewal Instructions</a>	2016-09-09
<a href="#">NIAGADS Data Use Certification</a>	2015-04-16
<a href="#">NIH Biosketch Sample</a>	2015-04-16
<a href="#">Research Use Statement Template</a>	2016-09-09
<a href="#">Supplemental Information</a>	2016-09-09
<a href="#">NIA AD Genomic Data Sharing Policy</a>	2015-03-02
<a href="#">Sample Language Data Return to NIAGADS</a>	2015-03-02

Data != Code

Data

Not protected by copyright

Data != Code

# Data

Not protected by copyright  
Patents do not apply to data

Data != Code

# Data

Not protected by copyright  
Patents do not apply to data  
Value in analysis

# Data != Code

## Data

Not protected by copyright  
Patents do not apply to data  
Value in analysis

## Code

Protected by copyright  
Standard patent process  
Value is intrinsic





# CDLA

Two model agreements, introduced and sponsored by the Linux Foundation

# CDLA

Two model agreements, introduced and sponsored by the Linux Foundation

October 23, 2017

Modeled after leading open source agreements

Licenses that reflect the nuances of data

Designed for use by independent data communities

CDLA promotes free exchange of data

## CDLA promotes free exchange of data

Permits data to be freely used, modified, and republished

Authorship and source attribution statements must be preserved

No use restrictions permitted

Broader license coverage than mere copyright

Explicit permissions to create separate works and analyses of licensed data

# Computational Use

Both CDLA licenses provide the right to Computational Use

# Computational Use

## Both CDLA licenses provide the right to Computational Use

- Analyze data and create analytical works based upon it

- Analytical works do not need to be relicensed under the terms of the CDLA

- Minimal reps and no warranties

- Broad Limitation of Liability

- No prohibition on commercial use of data

# CDLA Versions

Versions

Sharing

Permissive



## Versions

### Sharing

Republished data must be licensed under  
CDLA Sharing

### Permissive

Data may be republished under  
any terms not inconsistent with  
the License

## Versions

# Sharing

Republished data must be licensed under  
CDLA Sharing

### Includes:

- Data Modifications to data received
- Data Additions to data received

### Excludes:

- The results of any analysis

# Permissive

Data may be republished under  
any terms not inconsistent with  
the License

# Example



# Example

CDLA



# Example



# Example

CDLA - Sharing



CDLA - Sharing





# Example

CDLA - Permissive



# Example, pt. 1



# Deep Learning Developers Dataset

# Deep Learning Frameworks

Allow for “easy” construction of neural networks

# Deep Learning Frameworks

## Allow for “easy” construction of neural networks

Abstract away lowest level detail

Standardize commonly used neural network concepts, from components to models

Provide tooling for deployment on various hardware configurations, including GPUs

# Deep Learning Frameworks



TensorFlow

# Python Notebooks

```
In [1]: 1 !pip install pygithub

Requirement already satisfied: pygithub in /anaconda3/lib/python3.6/site-packages (1.39)
Requirement already satisfied: pyjwt in /anaconda3/lib/python3.6/site-packages (from pygithub) (1.6.1)
You are using pip version 18.1, however version 19.0.3 is available.
You should consider upgrading via the 'pip install --upgrade pip' command.

In [2]: 1 import re
2 import sys
3 import json
4 import base64
5 import github
6 import operator
7 import numpy as np
8 from collections import *
9 from github import GithubException

In [3]: 1 login = '***'
2 password = '***'
3
4 g = github.Github(login,password)

In [4]: 1 targetfiles = ['.ipynb']

In [5]: 1 def download_directory(repository, path, framework):
2     global dataset
3     try:
4         contents = repository.get_contents(path)
5         for content in contents:
6             if content.type == 'dir':
7                 download_directory(repository, content.path, framework)
8             else:
9                 if content.content:
10                    if len(str(content.name).split(".")) == 2:
11                        if any(substring == ("." + str(content.name).split(".")[1]) for substring in targetlangs):
12                            try:
13                                dataset.append([repository, str(base64.b64decode(content.content), 'utf-8'), csp])
14                            except (GithubException, IOError) as exc:
15                                print('Error processing %s: %s', content.path, exc)
16                    except (GithubException, IOError) as exc:
17                        print("error in dir ")

In [6]: 1 def getrepos():
2     repos = []
3     repos.append(list(g.search_repositories('TensorFlow')))
4     repos.append(list(g.search_repositories('PyTorch')))
```

# Notebook Licensing

Licenses vary widely among notebooks posted to GitHub

located in a separate file or repo

referenced in the notebook

no, ambiguous or default license (Cloud Platform)

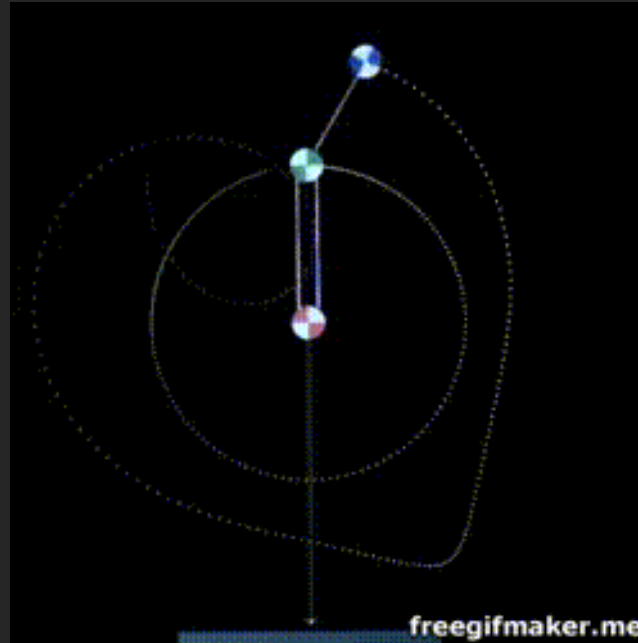
notebooks with one license referencing information with differing license

# Example, pt. 2

# Double Pendulum Chaotic Dataset



# Double Pendulum Chaotic Dataset



# Learning beyond simulated physics

## NIPS 2018 Spatiotemporal Workshop

Used to showcase ML – Prometheus Framework

Promotes reproducibility in Machine Learning Research

# The End

# Interested in data sharing?

# Thank You

Nick Acosta  
Developer Advocate  
—  
[nacosta@us.ibm.com](mailto:nacosta@us.ibm.com)  
[cdla.io](http://cdla.io)



# OPEN SOURCE

## LEADERSHIP SUMMIT

