Introduction
○○○○○○○○

History and Achievements
○○○○○○○○○○

Future
○○○○○○

# Machine Learning Open Source Software
## *(A PASCAL success story?)*

Sören Sonnenburg

Fraunhofer FIRST.IDA, Berlin

joint work with
*Cheng Soon Ong and Mikio Braun*

**FIRST**

**Fraunhofer** Institut
Rechnerarchitektur
und Softwaretechnik

Introduction
00000000

History and Achievements
0000000000

Future
000000

# Outline

Fraunhofer Institut
Rechnerarchitektur
und Softwaretechnik

# What is Open Source Software?

**Idea: Freedom to read, modify and redistribute source code**

MS Windows network stack, MacOSX (BSD based)

TV (Sharp HDTV Aquos)

Mobile Phones (Motorola RAZR, Android)

Wireless routers (Linksys WRT)

**Common: Free exchange of information, to avoid "reinventing the wheel".**

Fraunhofer Institut
Rechnerarchitektur
und Softwaretechnik

Introduction
○●○○○○○○○

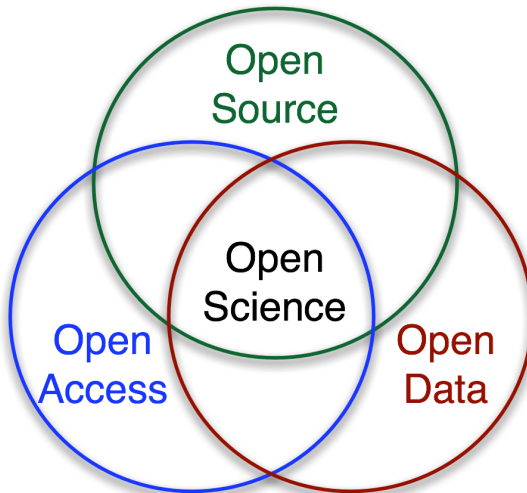History and Achievements
○○○○○○○○○○

Future
○○○○○○

Open Source?

# Open Source Definition (`www.opensource.org`)

- The Open Source Initiative (OSI) manages a license list of currently 65 approved open source licenses

**Criteria to be open source:**

1. Free redistribution
2. Must include source code
3. Derived works allowed
4. Integrity of the author's source code
5. No discrimination against persons or groups
6. No discrimination against fields of endeavor
7. License is redistributed
8. License must not be specific to a product
9. License must not restrict other software
10. License must be technology-neutral

Fraunhofer Institut
Rechnerarchitektur
und Softwaretechnik

**Introduction**
○○●○○○○○

History and Achievements
○○○○○○○○○○

Future
○○○○○○

Open Source and Science

# Open Science

# Open Access

> *Open access truly expands shared knowledge across scientific fields, it is the best path for accelerating multi-disciplinary breakthroughs in research.*
> *— Open letter to the U.S. Congress, signed by 25 Nobel laureates, (August 26, 2004)*

- Enabled by low-cost distribution on the Internet
- Open access literature is digital, online, free of charge, and free of most copyright and licensing restrictions. For example Creative Commons (creativecommons.org)
- Many journals (3096 according to www.doaj.org) have adopted the open access model (including JMLR, ...)

Fraunhofer Institut
Rechnerarchitektur
und Softwaretechnik

**Introduction**
○○○○●○○○

History and Achievements
○○○○○○○○○○

Future
○○○○○○

Open Source and Science

# Opening Machine Learning

- Hope to have a similar boost by adopting "open practices" in machine learning
  - *software and data* accompany paper
  - all openly licensed
- Some collections exists (UCI, Delve, Caltech, IDA Repository)
- How many machine learners publish software and data with their paper?
- Reasons? Misconception that open source renders commercial exploitation impossible?

**Focus on Machine Learning Open Source Software**

**Introduction**
○○○○○●○○

History and Achievements
○○○○○○○○○○

Future
○○○○○○

Open Source and Science

# Advantages of Machine Learning Open Source Software

- **Reproducibility** of scientific results
- **Fair comparison** of algorithms
- **Problems uncovered quickly**
- Building on existing resources (rather than re-implementing)
- **Access to scientific** tools without restrictions
- Easier to combine different advances
- **Faster adoption** of ML methods in other disciplines and in industry
- Collaborative emergence of standards

Fraunhofer Institut
Rechnerarchitektur
und Softwaretechnik

**Introduction**
○○○○○○○●○

History and Achievements
○○○○○○○○○○

Future
○○○○○○

## Obstacles to an MLOSS community

- Publishing software is **not considered a Scientific contribution**
- **Misconception** — Opening the source **conflicts with commercial interests**
- The **incentive** for publishing open source software is **not high enough**
- Machine learning researchers may **not be good programmers**
- **Sloppiness hides problems** of newly proposed methods and eases acceptance at conferences and journals.
- **Tradition** — reviewers pass papers of similar quality

Fraunhofer Institut
Rechnerarchitektur
und Softwaretechnik

Introduction
○○○○○○○●○

History and Achievements
○○○○○○○○○○

Future
○○○○○○

Obstacles

# Read our Position Paper

### The Need for Open Source Software in Machine Learning

**Sören Sonnenburg**[*]                     SOEREN.SONNENBURG@FIRST.FRAUNHOFER.DE
*Fraunhofer Institute FIRST*
*Kekulestr. 7*
*12489 Berlin, Germany*

**Mikio L. Braun**[*]                        MIKIO@CS.TU-BERLIN.DE
*Technical University Berlin*
*Franklinstr. 28/29*
*10587 Berlin, Germany*

**Cheng Soon Ong**[*]                        CHENGSOON.ONG@TUEBINGEN.MPG.DE
*Friedrich Miescher Laboratory*
*Max Planck Society*
*Spemannstr. 39*
*72076 Tübingen, Germany*

. . . Bengio, Bottou, Holmes, LeCun, Müller, Pereira, Rasmussen, Rätsch, Schölkopf, Smola, Vincent, Weston, Williamson

Introduction
○○○○○○○○

History and Achievements
●○○○○○○○○○○

Future
○○○○○○

History

# Timeline



Machine Learning Tools Satellite Workshop

Dec 2005       Dec 2006       July   Oct   Dec 2007

Introduction
○○○○○○○○

History and Achievements
○●○○○○○○○○

Future
○○○○○○

History

# Timeline

Introduction
○○○○○○○○○

History and Achievements
○○●○○○○○○○○

Future
○○○○○○

History

# Timeline

Introduction
○○○○○○○○

History and Achievements
○○○●○○○○○○

Future
○○○○○○

History

# Timeline



JMLR **Machine Learning Open Source Software**

Journal of Machine Learning Research 8 (2007) 2443-2466          Submitted 7/07; Published 10/07

**The Need for Open Source Software in Machine Learning**

| Dec | Dec | July | Oct | Dec |
| 2005 | 2006 | | | 2007 |

Introduction
○○○○○○○○

History and Achievements
○○○○●○○○○○

Future
○○○○○○○

History

# Timeline

Introduction
○○○○○○○○

History and Achievements
○○○○○○●○○○○

Future
○○○○○○

Achievements

# Workshops



Machine Learning Tools Satellite Workshop

First PASCAL workshop

- motivated by "The Mathworks" changing licenses. Affected people from Fraunhofer, Max-Planck, NICTA, INSA discussed and presented alternatives (octave, R, python,...)

Workshop on Machine Learning Open Source Software


Neural Information
Processing Systems
Conference

Second PASCAL/NIPS workshop

- Open call for papers. Received 20 submissions, 8 accepted. 3 invited speakers (Weka, scipy, cvxopt)
- Lively discussion, with the common themes:
  - incentives for researchers are missing
  - we should have a place to publish MLOSS
  - we still have a long way to go

Fraunhofer Institut
Rechnerarchitektur
und Softwaretechnik

Introduction
○○○○○○○○

History and Achievements
○○○○○○●○○○

Future
○○○○○○

Achievements

# New JMLR Track

**JMLR** **Machine Learning Open Source Software**

Contributions to http://jmlr.org/mloss/ should be related to

- Implementations of machine learning algorithms,
- Toolboxes,
- Languages for scientific computing

and should include

- A 4 page description,
- The code,
- A recognised open source license.

Fraunhofer Institut
Rechnerarchitektur
und Softwaretechnik

# Community site `mloss.org`

**All projects welcome**

- Implementations of machine learning algorithms,
- Toolboxes,
- Languages for scientific computing
- Data readers, preprocessing
- Concrete applications

and should include

- A recognised open source license.
- Pointer to project homepage and download link

**Contribute to `http://mloss.org`!**

Fraunhofer Institut
Rechnerarchitektur
und Softwaretechnik

Introduction
○○○○○○○○

History and Achievements
○○○○○○●○○●○

Future
○○○○○○

Achievements

# `mloss.org` Screenshot

Introduction
○○○○○○○○

History and Achievements
○○○○○○○○○●

Future
○○○○○○

Achievements

# The story so far...

Introduction
ooooooooo

History and Achievements
oooooooooo

Future
●ooooo

Success Story?

# Success Story? ⇒ Not yet

- JMLR track received 9 submissions, none accepted yet
  - Bigger established toolboxes are already published (e.g. Weka book)
  - It takes time to polish software projects to satisfy reviewers
  - New JMLR track not well known

- `mloss.org` currently has 180 registered users and 58 software projects
  - Collecting projects that are already out there
  - No collaborations / re-use of code yet
  - No lively discussion
  - Users mostly inactive

⇒ **How can we attract more (active) users?**

Fraunhofer Institut
Rechnerarchitektur
und Softwaretechnik

Introduction
○○○○○○○○

History and Achievements
○○○○○○○○○○

**Future**
○●○○○○

Success Story?

# Call for help

`mloss.org` **users needed**

- Use software
- Rate software
- Comment on software
- Discuss in the forum how we can improve
- Discuss about data standards etc, etc.

**Developers needed**

- Help us to implement data standards
- Submit your software to `mloss.org`
- Help us to further develop and maintain the website.

Fraunhofer Institut
Rechnerarchitektur
und Softwaretechnik

Introduction
○○○○○○○○

History and Achievements
○○○○○○○○○○

Future
○○●○○○

Success Story?

## Join the team!

**We are open and need**

- **Your Criticisms**
- **Your Ideas**
- **Your Feedback**
- **Your Contributions**

Fraunhofer Institut
Rechnerarchitektur
und Softwaretechnik

Introduction
00000000

History and Achievements
0000000000

Future
000●00

Success Story?

# Plans for 2008

**Events:**

- See publications in JMLR-MLOSS
- Get lively discussions in `mloss.org`
- NIPS'08 workshop (if it gets accepted, otherwise NIPS satellite workshop)

**Interoperability**

- A common data exchange format

Fraunhofer Institut
Rechnerarchitektur
und Softwaretechnik

# Data Set Standards

**Goal:** Develop a data exchange standard

- Currently, many data formats exist
    - ARFF
    - orange tab delimited
    - SVMlight, libsvm format
    - pyML, UCI
    - . . .
    - Post your thoughts!
      http://mloss.org/community/standards/13/
- A lot of time is wasted on converting data.

**First Proposal:** Use ARFF for dense vectorial data

- Used in Weka, code exists for R, matlab
- Subject to what the community thinks

Fraunhofer Institut
Rechnerarchitektur
und Softwaretechnik

Introduction
00000000

History and Achievements
0000000000

Future
00000●

# Summary

**Achievements**

- Organized two PASCAL workshops on MLOSS
- Established a JMLR track for MLOSS (so far 9 submissions)
- Position paper on "The need for Open Source Software in Machine Learning"
- Community site mloss.org (58 projects, 180 users)

**Future**

- Data exchange standards
- Shall we (How can we) extend this approach to Open Data ?
- New ideas; Where should we go beyond 2008?

Fraunhofer Institut
Rechnerarchitektur
und Softwaretechnik