# Building a scientific workbench in Pharo

Konrad Hinsen [1, 2]    Serge Stinckwich [3, 4, 5]

[1]Centre de Biophysique Moléculaire, Orléans, France

[2]Synchrotron SOLEIL, Saint Aubin, France

[3]Sorbonne Université, IRD, Unité de Modélisation Mathématiques et Informatique des Systèmes Complexes, UMMISCO, Bondy, France

[4]Université de Yaoundé I, Yaoundé, Cameroon

[5]Université de Caen Normandie, Caen, France

27 August 2019

# What's a scientific workbench?

- The IDE of the computational scientist
- Supports the tasks of doing science on a computer:
  - Write and test code
  - Import and export data
  - Process data
  - Perform simulations
  - Inspect experimental and computed data
  - Document all of the above
- Makes computations reproducible.

# RStudio

# Jupyter

localhost:8888/notebooks/influenza-like-illness-analysis-jupyter.ipynb

jupyter influenza-like-illness-analysis-jupyter Last Checkpoint: Last Monday at 5:03 PM (autosaved)  Logout

File   Edit   View   Insert   Cell   Kernel   Widgets   Help                          Not Trusted   Python 3 ○

First, we define the observation periods as the new index of our dataset. That turns it into a time series, which will be convenient later on.

Second, we sort the points chronologically.

In [7]:
```
sorted_data = data.set_index('period').sort_index()
```

We check the consistency of the data. Between the end of a period and the beginning of the next one, the difference should be zero, or very small. We tolerate an error of one second.

This is OK except for one pair of consecutive periods between which a whole week is missing.

We recognize the dates: it's the week without observations that we have deleted earlier!
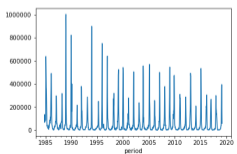
In [8]:
```
periods = sorted_data.index
for p1, p2 in zip(periods[:-1], periods[1:]):
    delta = p2.to_timestamp() - p1.end_time
    if delta > pd.Timedelta('1s'):
        print(p1, p2)
```

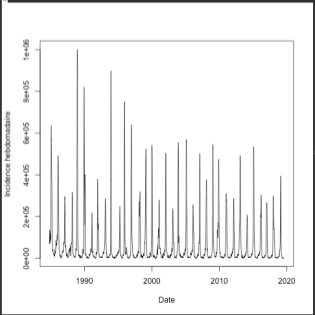1989-05-01/1989-05-07 1989-05-15/1989-05-21

A first look at the data!

In [9]:
```
sorted_data['inc'].plot()
```

Out[9]: <matplotlib.axes._subplots.AxesSubplot at 0x110f4b9b0>

# Emacs

# Shifting priorities

**Traditional focus:** *get work done efficiently*

- interactive computation
- generate plots and tables (for pasting into publications produced outside of the workbench)

**More recent criteria:** *robust and understandable results*

- reproducible computations
- shared/publishable raw datasets
- well-documented computations
- document while you compute

# The state of the art: computational notebooks

A fusion of scripts, REPLs, and literate programming, invented in the 1980's by Mathematica

- A linear sequence of so-called "cells"
- Three cell types:
    - Text cells hold rich text for documentation
    - Code cells contain code snippets
    - Output cells show the output of one code snippet (text or graphics)
- Code cells can be executed one by one, manually...
- ... or sequentially as part of a whole-notebook execution.

Many implementations: Mathematica, Jupyter, R Markdown, Emacs/Org-Mode, ...

# Limitations of notebooks

- Linear sequence of cells: no way to structure or modularize
- Made worse by shared mutable state...
- ... and even worse by interactive cell execution.
- Documentation follows code structure: no way to relegate technical details to an appendix
- Data dependencies are not explicit, nor easily visible.
- Neither code nor data are reusable by other notebooks.
- Different tools/user interfaces for notebooks and library code.

Notebooks blissfully ignore decades of software engineering achievements.

# Smalltalk to the rescue

Hypothesis:
    A Smalltalk system is a much better starting point
    for designing a scientific workbench than a REPL.

Nice properties:

- well-known to this audience!

Missing pieces:

- A documentation tool that allows embedding code and data.
- Management of computational tasks and data dependencies
  to replace the notebook's linear control flow
- Support libraries for scientific computing.

## SMALLTALK – THE NEXT GENERATION SCIENTIFIC COMPUTING INTERFACE?

Richard L. PESKIN, Sandra S. WALTHER and Andy M. FRONCIONI

*Center for Computer Aids For Industrial Productivity, Rutgers University \*, P.O. Box 1390, Piscataway, New Jersey 08855-1390, U.S.A.*

The need for rapid prototyping of numerical simulations is considered, and an object-oriented, graphical based system (Smalltalk) is proposed as a basis for a new approach to user interfaces for scientific computing. The interface system requirements for problem expression, automatic programming, visualization, computational steering, and concurrent computing are discussed.

## 1. Introduction

While scientific and engineering computation needs have been a major driving force in the

# Glamorous Toolkit

Nice properties:

- presented yesterday to this audience
- specifically for a scientific workbench: an excellent documentation tool

Missing pieces:

- Management of computational tasks and data dependencies
  to replace the notebook's linear control flow

# ActivePapers

A research project about performing and communicating computer-aided research

- Started in 2011.
- Initial focus: reproducible high-performance computing.
- Management of computational tasks and data dependencies
- Current implementation based on Python...
- ... and a lousy user interface: very basic CLI

# PolyMath

- Release 1.0 last week
- 300 classes, 50 packages, 24K LOC, 806 unit tests
- Ordinary differential Equations, Random Number Generators, Linear algebra, Matrices, Complex Numbers, FFT, Polynomials, Probability distributions, ...
- more recently: Automatic differentiation, Principal Component Analysis, t-SNE,
- DataFrame to do data analysis
- Talk on PolyMath next thursday

Pharo + Glamorous Toolkit + PolyMath + ActivePapers
= Scientific Workbench

# Outlook

A lot of work remains to be done:

- More domain-specific libraries
- Interfaces to other languages
- Data management outside of the Pharo image
- Publishing ActivePapers on the Web