# Modeling spontaneous brain activity in Python
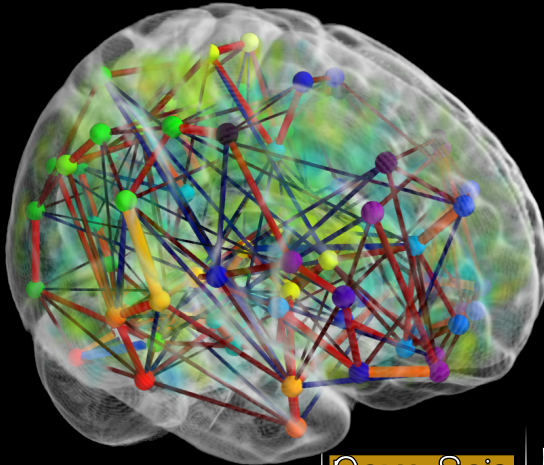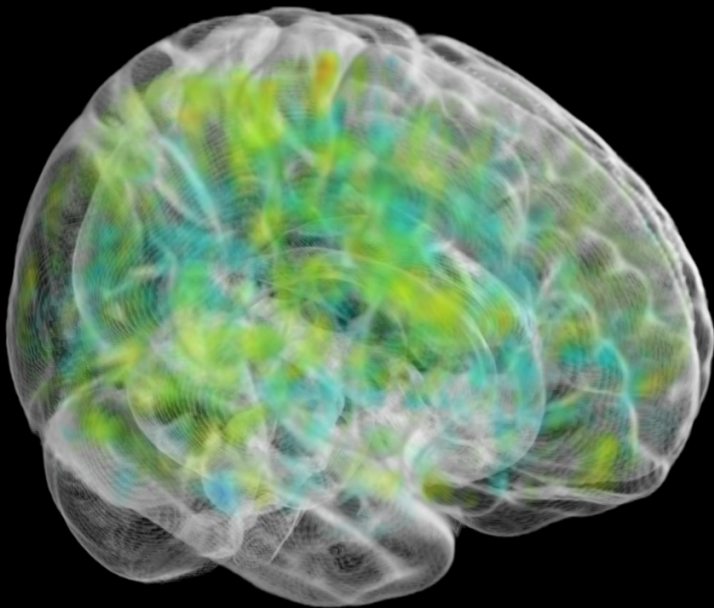## Scientific progress and software challenges

**Gaël Varoquaux**, INRIA and Neurospin
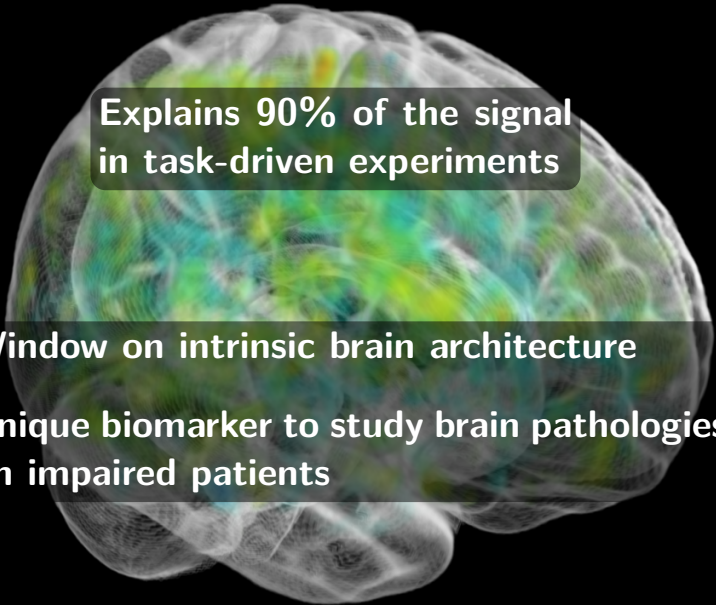
# Why study spontaneous brain activity?

Explains 90% of the signal in task-driven experiments

- Window on intrinsic brain architecture

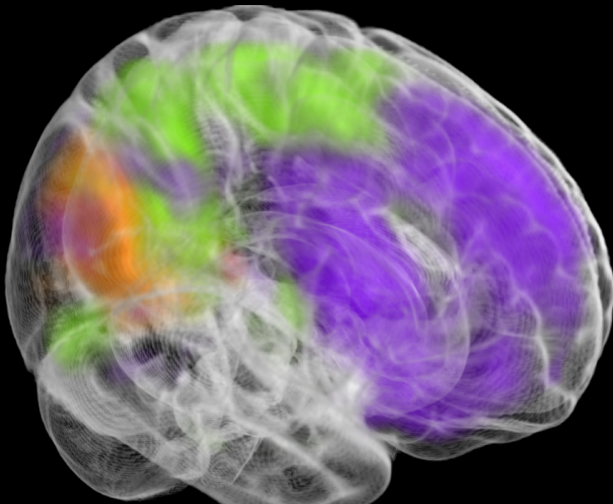- Unique biomarker to study brain pathologies on impaired patients

**Scientific challenge**

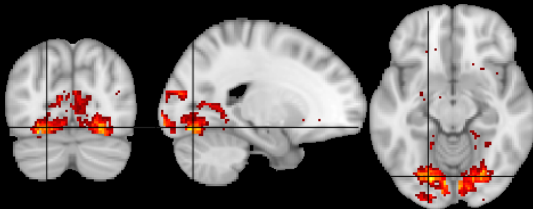To develop in collaboration with neuroscientists new statistical tools to learn probabilistic models of spontaneous brain activity

■ Study of stimuli response



■ Mass-univariate statistics:
  for each voxel   $\mathbf{X} = \beta \mathbf{Y} + \mathbf{E}$

■ Group inference: subject-variability model on $\beta$

**Nipy**: NeuroImaging in Python
Berkeley, Stanford, Neurospin ...
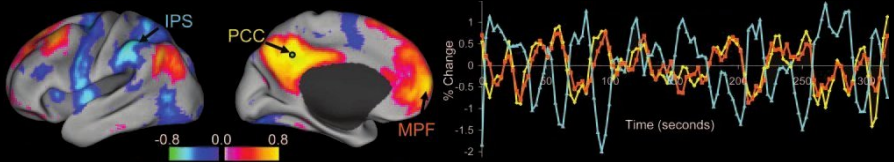
Vision: Open code shared between labs

Progress: ■ Statistical models implemented 🙂
■ API difficult to use 😫
■ Good Input/Output code 🙂
■ Preprocessing not implemented 😫

Roadblocks: ■ Different teams ⇒ different visions
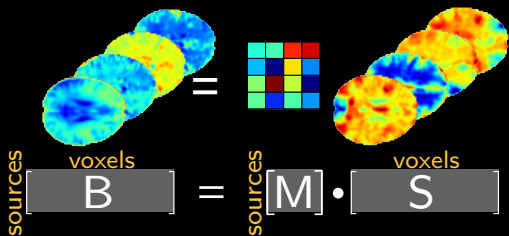■ Scientists can't justify time on "solved problems"

- Biswal 1995: strong correlation between activity in left and right motor cortex at rest

- Later: seed-based correlation mapping



*The human brain is intrinsically organized into dynamic, anticorrelated functional networks* (Fox 2005)
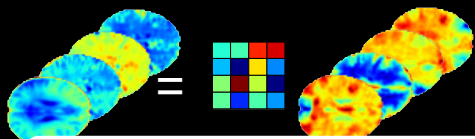
**How many? How to choose seeds?** 😣

**B**: observed images

**M**: mixing matrix

**S**: sources

$$B = M \cdot S$$

- Minimize mutual information between patterns $S$.

**B**: observed images

**M**: mixing matrix

sources

$S$.

**B**: observed images

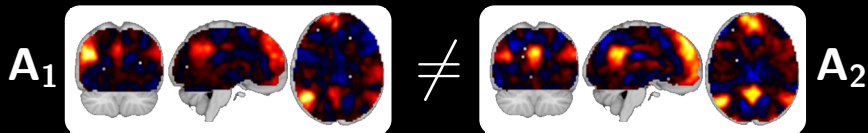**M**: mixing matrix

**S**: sources

$$B = M \cdot S$$

- Minimize mutual information between patterns $S$.

## No noise model

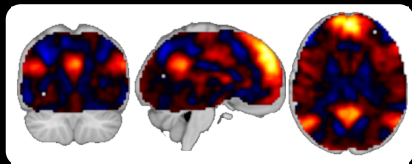$\Rightarrow$ **Lack of reproducibility** $+$ **Fits noise**

# 1 Model subject-to-subject variability

$A_1$ $\neq$ $A_2$

Multivariate random effects model:

$$\mathbf{Y_s} = \text{loadings} \times \mathbf{P_s} + \text{intra-subject noise} \qquad \text{PCA}$$
$$\{\mathbf{P_s}\} = \text{loadings} \times \mathbf{B} + \text{inter-subject variability} \qquad \text{CCA}$$
$$\mathbf{B} = \mathbf{M} \times \mathbf{A} \qquad \text{ICA}$$

$\Rightarrow$ **Group-level networks**

# 1 Model subject-to-subject variability

**Reproducibility across random groups**

|            | no CCA     | CCA + ICA  |
|------------|------------|------------|
| Subspace   | .36 (.02)  | .71 (.01)  |
| One-to-one | .36 (.02)  | .72 (.05)  |

*Varoquaux*, NeuroImage 2010

**Problem to solve**:

(1)     $\mathbf{Y_s} = \text{loadings} \times \mathbf{P_s} + \ldots$    PCA: SVD

(2)   $\{\mathbf{P_s}\} = \text{loadings} \times \mathbf{B} + \ldots$    CCA: SVD

(3)     $\mathbf{B} = \mathbf{M} \times \mathbf{A}$        ICA: iterations

$+$ Recomputed many times across random groups

---

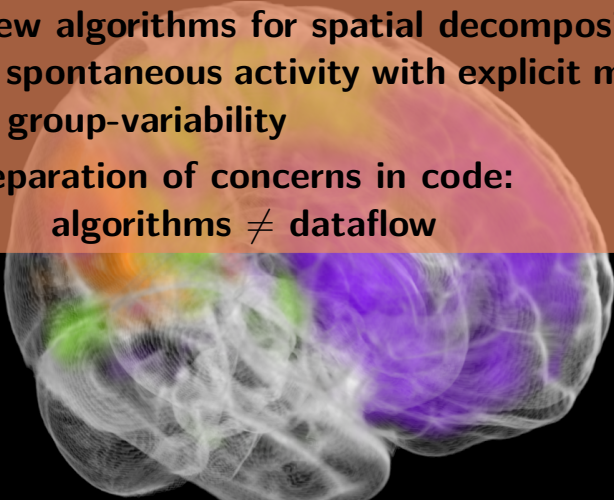Step 2 and 3: Small data size $\Rightarrow$ not bottleneck

Step 1: ■Independent problems per subject
$\Rightarrow$Parallel runs and caching of the results
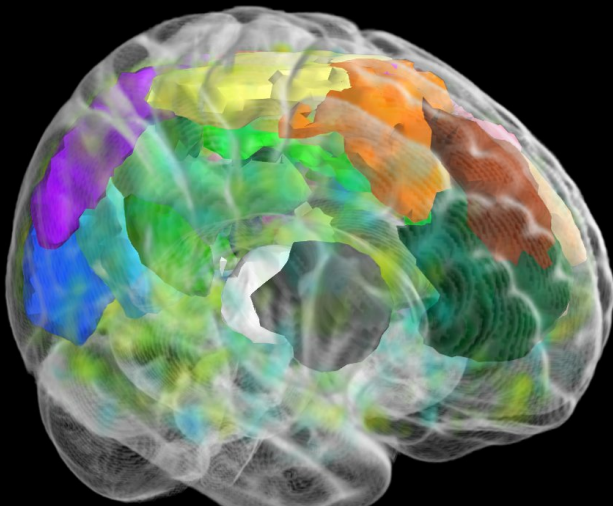
**Joblib:** **Python functions as pipeline jobs**

**Goals: remove dataflow and persistence problems from algorithmic code**

## Spatial patterns of brain activity

- New algorithms for spatial decomposition of spontaneous activity with explicit model of group-variability
- Separation of concerns in code:
  algorithms $\neq$ dataflow

**2** **Beyond activation maps**

sources $[\quad B \quad] = $ sources $[M] \cdot [\quad S \quad]$

voxels · · · voxels

$$\text{sources}\left[\begin{array}{c}\overset{\text{voxels}}{B}\end{array}\right] = \text{sources}\left[M\right]\bullet\left[\begin{array}{c}\overset{\text{voxels}}{S}\end{array}\right]$$

$$\text{sources}\left[\begin{array}{c}\overset{\text{voxels}}{B}\end{array}\right] = \text{sources}\left[M\right]\bullet\left(\left[\begin{array}{c}\overset{\text{voxels}}{S}\end{array}\right] + \left[\begin{array}{c}\overset{\text{voxels}}{Q}\end{array}\right]\right)$$

- Interesting sources $S$ **sparse**
- $Q$: Gaussian noise

  $\Rightarrow$ Null hypothesis: centered normal distribution.

## Visual system



map 0, reproducibility: 0.54

-74    0    9

**V1**

map 1, reproducibility: 0.52

-91    3    -3

**V1-V2**

map 3, reproducibility: 0.47

-80    40    4

**extrastriate**

map 25, reproducibility: 0.34

-78    10    24

**superior parietal**

## Motor system



map 4, reproducibility: 0.47

part of motor

-25          -1        62

map 21, reproducibility: 0.36

part of motor

-21          -42        54

map 32, reproducibility: 0.30

part of motor

-8          -54        29
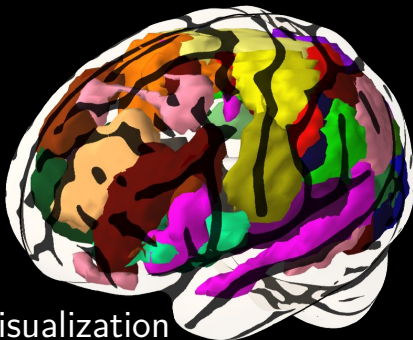
## Correlation matrix Σ
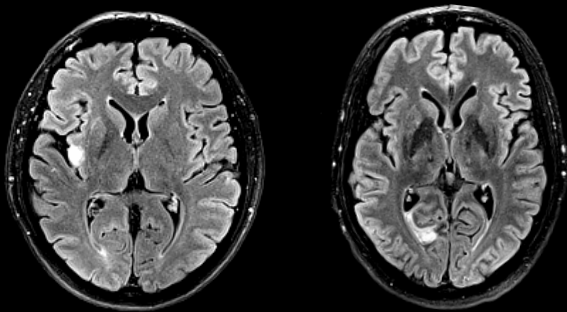
Data

Visualization

## Change of representation

**Understanding complex data requires inter-active visualization with** *high level concepts*
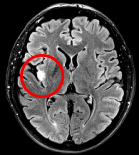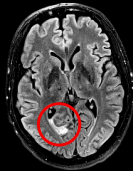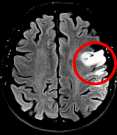
**Mayavi:**
Python 3D visualization

**Ischemic stroke**: 
- Temporary interruption of blood flow
- Affects 1 person out of 100 every year for people > 55 years
- Causes focal lesions of varying consequences

motor deficiencies    language impairments    coma    . . .



**How does brain reorganize after stroke?**

**Prognostic based on intrinsic brain activity?**

**Probabilistic model of data**

- Covariance = $2^{nd}$ moment of observed data

$\Rightarrow$ Specifies a probability distribution

Test the likelihood of data in a covariance model

## Probabilistic model of data

- Covariance = $2^{nd}$ moment of observed data

$\Rightarrow$ Specifies a probability distribution

Test the likelihood of data in a covariance model

---

## Covariances variations in healthy population



Which one of the above has a large cortical lesion?

## Probabilistic model of data

■ Covariance = $2^{nd}$ moment of observed data

⇒ Specifies a probability distribution
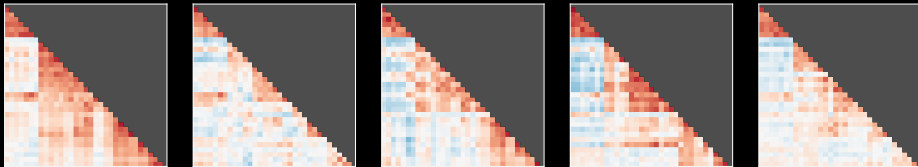
Test the likelihood of data in a covariance model

---

## Covariances variations in healthy population



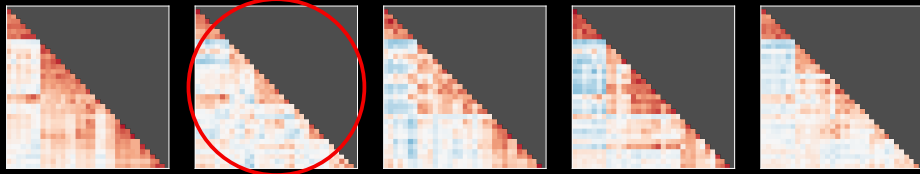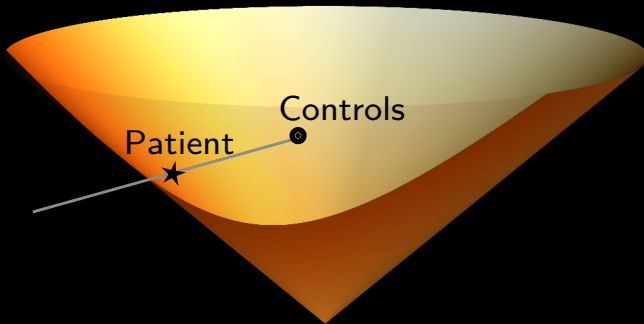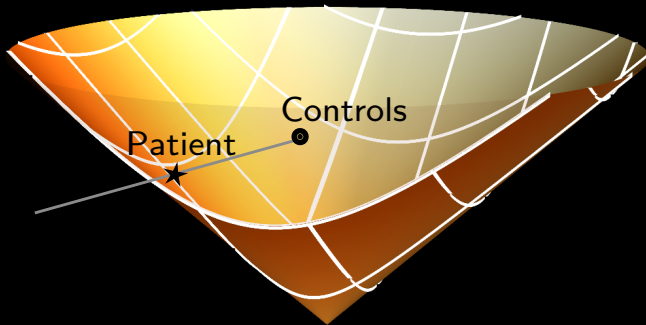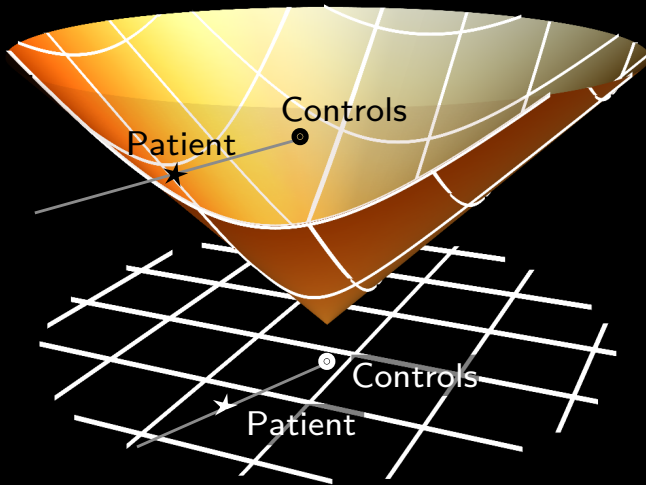Which one of the above has a large cortical lesion?

Controls

Patient

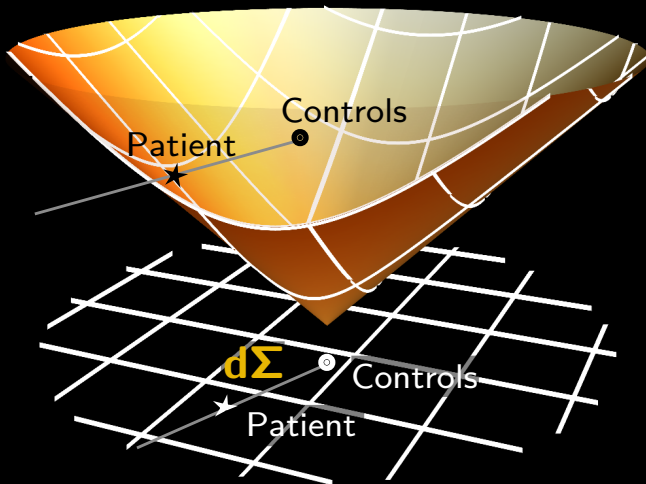Controls

Patient

Controls

Patient

Controls

Patient

$\mathcal{P}(\mathbf{d\Sigma})$: **probability density in tangent space**

$\mathcal{P}(\mathbf{d\Sigma})$: **probability density in tangent space**

## Between which regions is connectivity is modified?

### Ill-posed problem
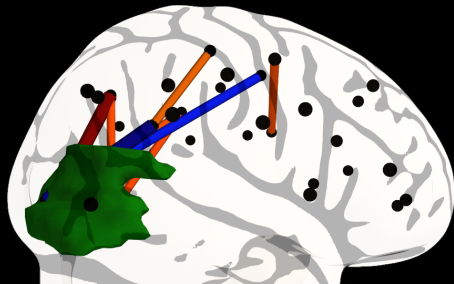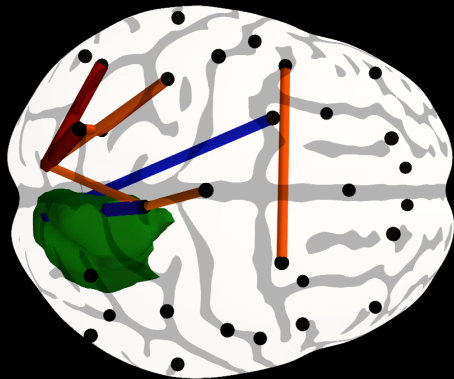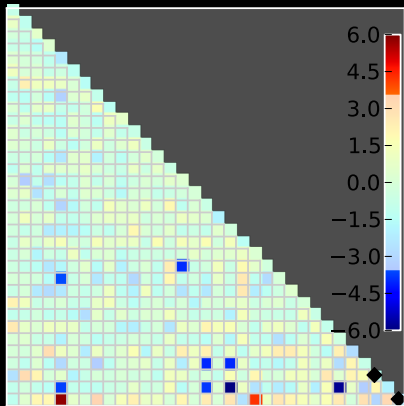■ Non-local effects

⇒ Many differences causes give the same observations

### Our suggestion
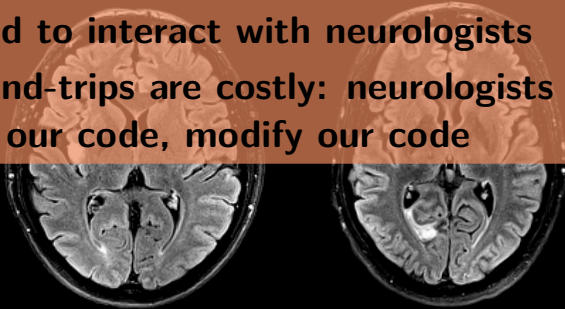■ Pair-wise partial correlations
■ In tangent space: almost independent
■ Draw random groups of healthy controls to tabulate their variability

# Research code in clinical settings

- Applications give rise to non-trivial mathematical problems
- Need to interact with neurologists
- Round-trips are costly: neurologists should use our code, modify our code

**4** **From models to software tools?**

- Gap from paper to software:
  Remove duplication          Write documentation
  Make usable APIs     Write tests     Fix corner cases

## Cost of code

Complexity scales as the square of project size

Woodfield 1979, *an experiment on unit increase in problem complexity*

## Cost of users

- Backward compatibility
- Support for multiple installations and versions
- Bug reports, feature request, mailing list support

$$\text{Maintenance cost} \sim (\text{\# lines})^2 \sqrt{\text{\# users}}$$

## Better code

- High-level coding and abstractions
  - numpy arrays: abstract out memory and pointers
  - traits Model+View: hide dialogs and events
  - joblib: factor out persistence
- Common libraries
  scipy, Mayavi, . . .

## Project management decisions

- 80/20 rule
- Not every research code should be released
- Focus on documentation and installation

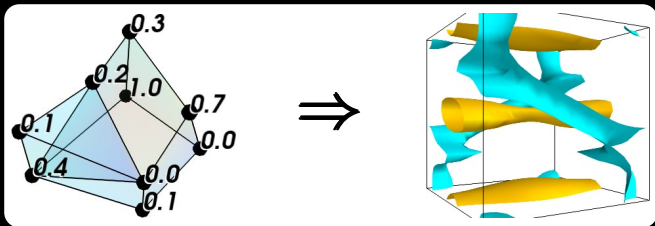## Segregated, functionally-specialized, packages

- Answer a specific problem
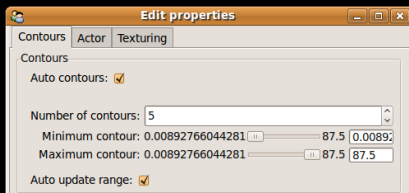- Limit dependencies

## Reusable projects

- Useful for a different purpose than the original one
- Libraries (no control of point of entry)
- Standard data structures
- Most often simple
- BSD licensed

# 4 Mayavi: making 3D visualization reusable

**Pipelines**: from data sources to visualization objects



■ Simple API: `mlab.contour3d(x, y, z, data)`

■ Building pipelines by function calls:
`mlab.pipeline.iso_surface(mlab.pipeline.contour(src))`

■ GUI
  + automatic script generation

**4** Mayavi: making 3D visualization reusable

260 lines of code!

**4 Mayavi: making 3D visualization reusable**

260 lines of code!

■ **All dialogs are components:** we expose our internals

■ **Visualizations included Traits view**

■ **Easy update of data**

Dataflow pipeline: *succession of processing steps executed on demand*

---

**joblib**:
- Lazy-revaluation
- Persitence
- Parallel processing
- Logging

**All with functions (seemingly)**

# 4 joblib: not writing pipelines

```
>>> from joblib import Memory
>>> mem = Memory(cachedir='/tmp/joblib')
>>> import numpy as np
>>> a = np.vander(np.arange(3))
>>> square = mem.cache(np.square)
>>> b = square(a)

_____
[Memory] Calling square...
square(array([[0, 0, 1],
       [1, 1, 1],
       [4, 2, 1]]))
_____square - 0.0s, 0.0min

>>> c = square(a)
>>> # The above call did not trigger an evaluation
```

# Towards Quantitative modeling of spontaneous brain activity

- Requires probabilistic models and state-the-art machine learning tools

- Algorithms and software development hand in hand with neurologists for applications

- Need a high-level stack of software tools general purpose with separation of concerns