# Open Questions for
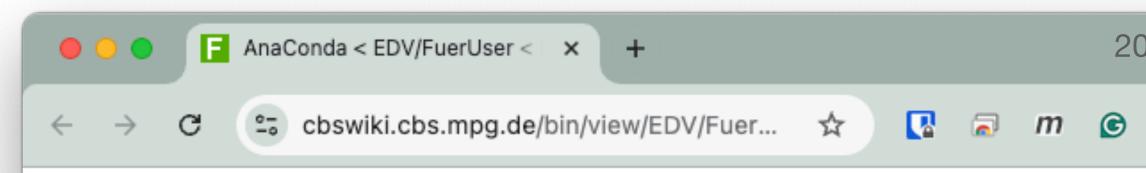# Computing Language for Open Science

**NCML lab meeting**

**2024-10-02 Seung-Goo KIM**

# Motivation
## Anaconda drama



**2024-09-03**

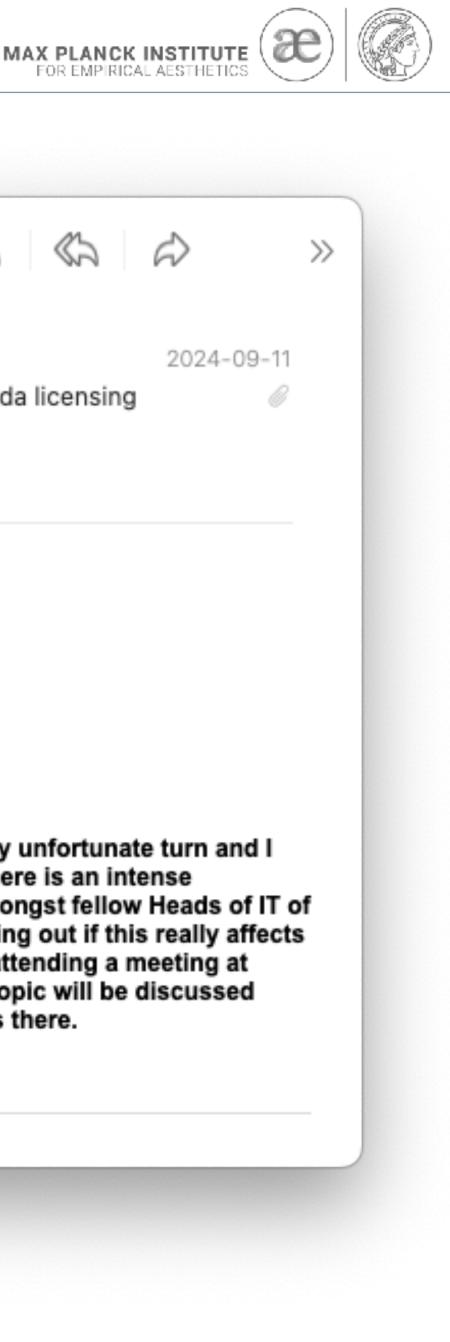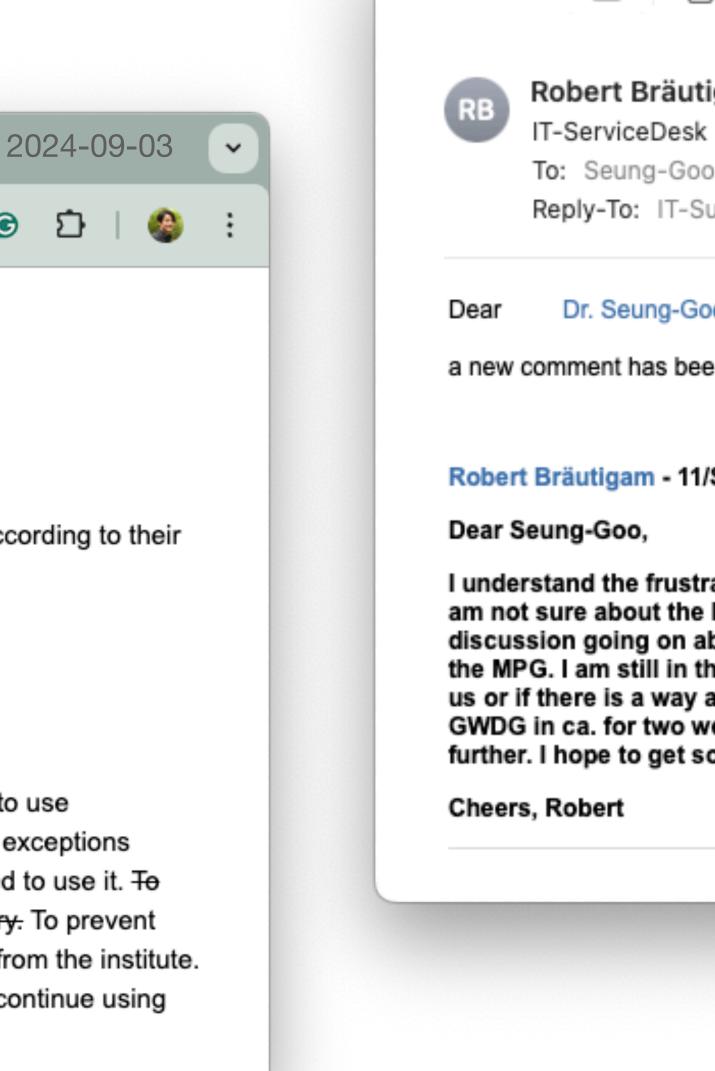cbswiki.cbs.mpg.de/bin/view/EDV/Fuer...

# ANACONDA

Permanent Link:

**Summary:** Anaconda is a Python distribution featuring the package manager *conda*. According to their terms of service it cannot be used at the institute without a license.
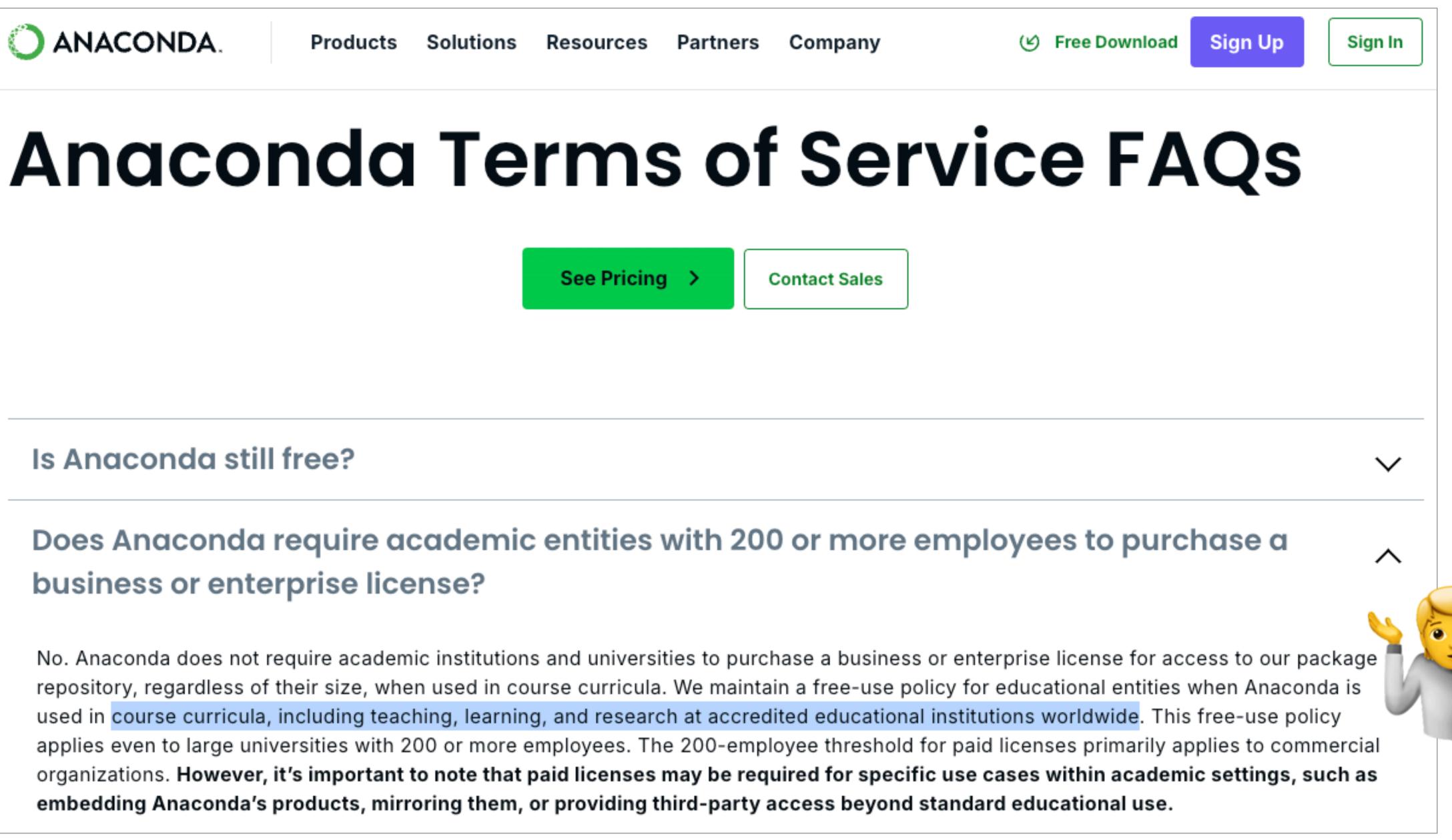
## Can I use Anaconda at the institute?

No.

According to their terms of service we (the Max Planck Society (MPS)) are not allowed to use "Anaconda's offerings" for free as we have more than 200 employees and none of their exceptions applies to our institute. Since the MPS has no licenses for Anaconda you are not allowed to use it. ~~To prevent license violations we are blocking connections to Anaconda's package repository.~~ To prevent disruption of the research process, the Anaconda package repository is still accessible from the institute. Please migrate away from channels at `repo.anaconda.com` as soon as possible! To continue using conda you can switch to Miniforge and the conda-forge Channel.



**Robert Bräutigam (Jira)**     2024-09-11
IT-ServiceDesk  IT-17358 Anaconda licensing
To:  Seung-Goo Kim,
Reply-To:  IT-Support

Dear     Dr. Seung-Goo Kim  ,

a new comment has been added:

**Robert Bräutigam** - **11/Sep/24 8:37 AM**

**Dear Seung-Goo,**

**I understand the frustration, this is a very unfortunate turn and I am not sure about the legal situation. There is an intense discussion going on about this topic amongst fellow Heads of IT of the MPG. I am still in the process of finding out if this really affects us or if there is a way around this. I am attending a meeting at GWDG in ca. for two weeks, where this topic will be discussed further. I hope to get some more insights there.**

**Cheers, Robert**

ANACONDA.

Products     Solutions     Resources     Partners     Company

Free Download     Sign Up     Sign In

# Anaconda Terms of Service FAQs

See Pricing  >      Contact Sales

## Is Anaconda still free?   ⌄

## Does Anaconda require academic entities with 200 or more employees to purchase a business or enterprise license?   ⌃

No. Anaconda does not require academic institutions and universities to purchase a business or enterprise license for access to our package repository, regardless of their size, when used in course curricula. We maintain a free-use policy for educational entities when Anaconda is used in course curricula, including teaching, learning, and research at accredited educational institutions worldwide. This free-use policy applies even to large universities with 200 or more employees. The 200-employee threshold for paid licenses primarily applies to commercial organizations. **However, it's important to note that paid licenses may be required for specific use cases within academic settings, such as embedding Anaconda's products, mirroring them, or providing third-party access beyond standard educational use.**

# What computing language should we use?

# **What computing language should we use** for open science?

# Why do I associate Python with Open Science?

- Until 2000s?: We used proprietary software all the time (MATLAB, BESA, BrainVoyager, SPSS, ..., Microsoft Office, Windows, Mac{HW+SW})

- 2012: Special Section on "Replicability in Psychological Science" in Perspectives on Psychological Science, https://journals.sagepub.com/doi/10.1177/1745691612465253

- 2015: Open Science Collaboration, "Estimating the reproducibility of psychological science", Science, https://doi.org/10.1126/science.aac4716

- For **transparency**, the whole analysis needs to be replicable on others' system only using open-access software (i.e., docking everything won't violate any copyright laws) => "MATLAB is proprietary🤑. Thus, MathWorks is the enemy of OpenScience👹!"

- Now if we need to pay for Anaconda (or maybe Python🐍 itself), it would make switching to Python meaningless for Open Science!

# Skepticism triggered! 🤔🧐🤨

## A punch in the face wakening up from the naiveness

- Is every open-access project👹 just a long-term investment for the market dominance?

- Is every public talk😈 an advertisement for a book, a promotion for a product, or a propaganda for a vote?

- "*If you're not paying for the product, then you are the product.*" (Tristan Harris🧌) (or is he just selling his books?)

- (Am I scamming you right now?😰🙀)

- ...but let's get real and let's think about the problem a bit more

# Why did I start using those languages?
## But first...

- **MATLAB**: because of SPM, SurfStat, EEGLAB, FieldTrip, and other popular toolboxes

- **Python**: because of MNE-python, TensorFlow, and other popular packages

- **R**: because of linear mixed-effects models (lme4)

- **Bash**: because of FSL, ANTs, FreeSurfer

- I didn't choose those language. I only chose packages.

- The developers of the popular packages/libraries chose their languages!

- So, how much meaningful is it for me to consider which language is good or bad?🤔

# Are we end-users or developers?

## How much helpful to learn languages?

- Some experimental physicists (like CERN or NASA) build their own tools [e.g., large hadron colliders, space telescopes, ...].

- We mainly develop theories, designs, and analyses, but also sometimes "design" new devices/environments [e.g., MR-piano, ArtLab].

- [my take]: Normally we use software as a tool, but sometimes it helps if we know how to tweak a little (it is also a bit risky though).

https://www.europosters.de/poster/rick-morty-portal-v35685

# Open questions

- Is proprietary software evil?

- Can the transparency be implemented only by open-access software?

- What is the best language?

# Open questions

- **Is proprietary software evil?**

- Can the transparency be implemented only by open-access software?

- What is the best language?

# Is capitalism inherently evil?
## Why some people feel software should be ethically free?

*Those who have plenty want more and so lose all they have.*

St. IGNUcius, the Church of EMACS (2012)
Free Software Foundation (1985-)

The Goose that Laid the Golden Eggs
Aesop's fable (~600 BCE), wiki

Bill Gates with red eyes, r/linuxmemes

# How do they make open access software?

## Are they just heavenly creatures living off their own niceness?

- **Linus Torvalds** (the creator of Linux) gets 1.7M USD in compensation as a Fellow of Linux Foundation

- **Guido van Rossum** (the creator of Python) worked at own startups, and Google[05-12], Dropbox[13-19], Microsoft[20-], 1-5M USD/yr?

- **Ross Ihaka** & **Robert Gentleman** (the co-creators of R) worked as statistic professors at University of Auckland, NZ

- **Travis Oliphant** (the creator of NumPy, SciPy and a co-founder of NumFOCUS, Anaconda[CEO:12-17]) gets compensation from as a staff at NumFOCUS.

- **Bjarne Stroustrup** (the creator of C++) worked as a programmer at Bell Labs, as a professor at Texas A&M University, and Columbia University.

# When you see a nice open-access project...

## How long will it be maintained and last open-access?

- A non-profit foundation?

- A for-profit start-up?

- A fixed-term research funding?

- A tenured professor's hobby?

# Is software public goods or common goods?

- Excludability: a possibility to exclude access of non-owners

- Rivalry: one's usage does not affect the availability for others

- Goods:

|  | Excludable | Non-excludable |
|---|---|---|
| **Rivalrous** | Private goods | Common goods (e.g., open data?) |
| **Non-rivalrous** | Club goods (e.g., proprietary s/w) | Public goods (e.g., open source s/w) |

# Common goods

## In capitalistic economics, it's called "externality"

- Hardin G, 1968, The Tragedy of the Commons, *Science*.

  - Depletion of public resources is inevitable (private or public ownership is needed)

- Dietz et al., 2003, The Struggle to Govern the Commons, *Science*.

  - Can be avoided by communal efforts and other institutions (i.e., making them club goods)



The Tragedy of the Commons

Use of the commons is below the carrying capacity of the land. All users benefit.

If one or more users increase the use of the commons beyond its carrying capacity, the commons becomes degraded. The cost of the degradation is incurred by all users.

Unless environmental costs are accounted for and addressed in land use practices, eventually the land will be unable to support the activity.

https://www.sustainable-environment.org.uk/Earth/Commons.php

# Your thoughts? 🧐💭

# Open questions

- Is proprietary software evil? (should/can we be non-evil in the long run?)

- Can the transparency be implemented only by using open-access languages?

- What is the best language?

# MATLAB for Open Science?

**At least they also know it "sells" for now**

# CodeOcean: Open Science Library
## It's free for now... but why not institute-level license?

**For journals**



**For authors**

# Your thoughts? 🧐💭

# Open questions

- Is proprietary software evil? (should/can we be non-evil in the long run?)

- Can the transparency be implemented only by open-access software?

- What is the best language?

Years (left axis): 1956, 1958, 1960, 1962, 1964, 1966, 1968, 1970, 1972, 1974, 1976, 1978, 1980, 1982, 1984, 1986, 1988, 1990, 1992, 1994, 1996, 1998, 2000, 2002, 2004, 2006, 2008, 2010, 2012, 2014, 2016, 2018

Nodes: Fortran I, Lisp, COBOL, Algol 60, PL/I, Smalltalk, Pascal, Prolog, Scheme, ML, Fortran 77, C (K&R), Smalltalk 80, Ada 83, Common Lisp, C++, MATLAB, SML, Eiffel, Perl, Caml, Tcl, Python, Fortran 90, Ruby, Java, JavaScript, Perl 5, R, OCaml, Haskell 98, Scheme R5RS, C++ (ISO), Python 2.0, C#, C# 2.0, Java 5, Go, Haskell 2010, Rust, Kotlin, C# 5.0, Julia, Swift, Java 8, JavaScript ES2017

25

# What is the best computing language?
## Computing languages as meta-tools to build your own tools

- As an end-user of a language, we want tools with high...

  - usability [HCI]: easy to learn/read/write/debug/maintain

  - accessibility [OS]: can be used without paid licenses

  - sustainability [OS, Econ]: 50+ years

  - performance [CS]: can handle large data

- ... in a good balance.

| Language | Latest release | GPL/DSL | Usability | Accessibility | Sustainability | Performance |
|---|---|---|---|---|---|---|
| **Python** | 2024-10-01 | GPL | ? | (600$/yr/user?) | ? | mid |
| **GNU-c-LISP** | 2023-01-13 | GPL | ? | 0 | ? | ? |
| **C** | 2024-02-21 | GPL | ? | 0 | ? | high |
| **C++** | 2023-03-19 | GPL | ? | 0 | ? | high |
| **C# (Win)** | 2023-11-14 | GPL | ? | 0 | ? | high |
| **Swift (Mac)** | 2024-03-01 | GPL | ? | 0 | ? | high |
| **Julia** | 2024-10-01 | GPL | ? | 0 | ? | mid |
| **Java** | 2024-09-17 | GPL | ? | 180$/yr/user | ? | high |
| **R** | 2024-06-14 | DSL | ? | 0 | ? | ? |
| **GNU Octave** | 2024-06-07 | DSL | ? | 0 | ? | mid |
| **Fortran** | 2023-11-17 | DSL | ? | 0 | ? | high |
| **Mathematica** | 2024-07-01 | DSL | ? | 839$/yr/user | ? | ? |
| **Stata** | 2023-04-25 | DSL | ? | 925$/yr/user | ? | ? |
| **WolframOne** | 2024-07-31 | DSL | ? | 1710$/yr/user | ? | ? |
| **MATLAB** | 2024-09-12 | DSL | ? | 263-3k$/yr/user | ? | mid |

# HCI research
## ChatGPT4o says...

- The **Cognitive Dimensions of Notations (CDN)** is a widely-used framework for evaluating the usability of notational systems, including DSLs (domain-specific languages). It provides a structured way to analyze how the design of a DSL affects the cognitive load on users. Some key dimensions include:

  - **Closeness of Mapping**: How well the language matches the mental model of the domain.

  - **Consistency**: How predictable and uniform the language is.

  - **Abstraction Gradient**: The complexity and flexibility in expressing concepts.

  - **Error-proneness**: How likely users are to make mistakes.

  - **Viscosity**: Resistance to change or how hard it is to make modifications.

# Closeness of Mapping: Let's invert a matrix !

$$A^+, (a_{i,j}) \in \mathbb{R}^{n \times m}, n > m$$

| <Python> | <R> | <octave; MATLAB; Julia> |
|---|---|---|
| ```python`<br>import numpy as np`<br>`A = np.random.rand(10, 5)`<br>`A_pinv = np.linalg.pinv(A)````` | ```r`<br>A <- matrix(runif(50), 10, 5)`<br>`A_pinv <- MASS::ginv(A)````` | ```matlab`<br>A = rand(10, 5)`<br>`A_pinv = pinv(A)````` |

| <APL> | <Perl+PDL> | <Cpp> |
|---|---|---|
| ```apl`<br>A ← 10 5 ρ ? 100`<br>`+⌹A````` | ```perl`<br>use PDL;`<br>`$A = random 10,5;`<br>`$A_pinv = pinv($A);````` | ```cpp`<br>#include <iostream>`<br>`#include <armadillo>`<br>`using namespace arma;`<br>`int main() {`<br>`    mat A = randu<mat>(10, 5);`<br>`    mat A_pinv = pinv(A);````` |

| <Fortran> | <LaTex+SageMath> | <HTML+MathJS> |
|---|---|---|
| ```fortran`<br>program pseudoinverse`<br>`  implicit none`<br>`  integer, parameter :: m = 10, n = 5`<br>`  real(8), dimension(m,n) :: A`<br>`  real(8), dimension(n,n) :: A_pinv`<br>`  real(8), dimension(m) :: singular_values`<br>`  real(8), dimension(m,m) :: U`<br>`  real(8), dimension(n,n) :: VT`<br>`  real(8), dimension(n) :: S_inv`<br>`  integer :: info, i, j`<br>`  real(8), external :: random_number`<br>``<br>`  call random_seed()````` | ```latex`<br>\documentclass{article}`<br>`\usepackage{sagetex}`<br>`\begin{document}`<br>`\begin{sagesilent}`<br>`A = random_matrix(RDF, 10, 5)`<br>`A_pinv = A.pseudoinverse()`<br>`\end{sagesilent}`<br>`\[ \sage{A} \sage{A_[pinv} \]`<br>`\end{document}````` | ```html`<br><!DOCTYPE html>`<br>`    <script src="https://cdnjs.cloudflare.com/ajax/libs/mathjs/10.6.4/math.min.js"></script>`<br>`<body>`<br>`    <script>`<br>`      let A = math.random([10, 5]);`<br>`      let A_pinv = math.pinv(A);`<br>`document.getElementById('output').innerHTML = "Original Matrix (10x5):\n" + math.format(A, {precision: 3}) + "\n\nPseudoinverse:\n" + math.format(A_pinv, {precision: 3});`<br>`    </script>`<br>`</body>````` |

30

# My conclusions

- Still Python & R seem to be good options.

- Non-CS scientists love math-oriented languages like R, MATLAB, Julia (we still misunderstand that we are handling natural numbers)

- Julia looks very interesting (e.g., UTF-8 for all names: 🤨`(a,b) = a^b`; 파이=pi; 🤨`(파이,3)`; built-in functions for math operations like MATLAB)

- MATLAB is so expensive… but so as many other tools we use (like Macs and Adobe suites).

- Usability is difficult to measure; and the field-leading developers decide which language we use.

- Transparency may be still achieved using specialized platforms like CodeOcean (it's not a trivial task that many individual researchers can handle)