

Pas d'IA sans data.

Vers des jeux de données libres et ouverts pour l'éducation.

 @jmaupetit //  @open_fun

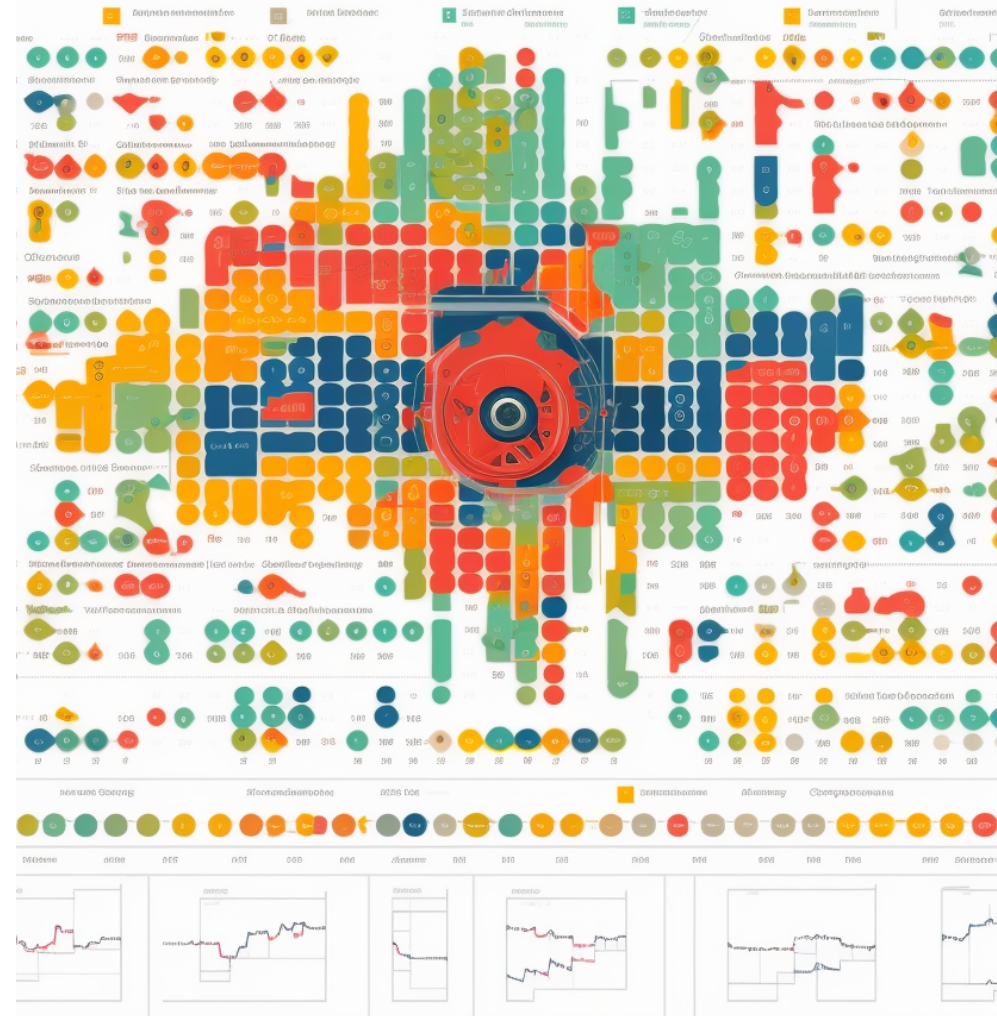


Cette présentation ne parle pas d'« Intelligence Artificielle ».

De l'importance des jeux de données pour apprendre

Problématique

- Dépendance : un jeu de données comportementales
- Comment obtenir un tel jeu de données ?
 - données maison
 - données ouvertes
- Très peu de jeux de données sont disponibles (OULAD, MOOCCubeX, StanfordMoocPosts, et ?)



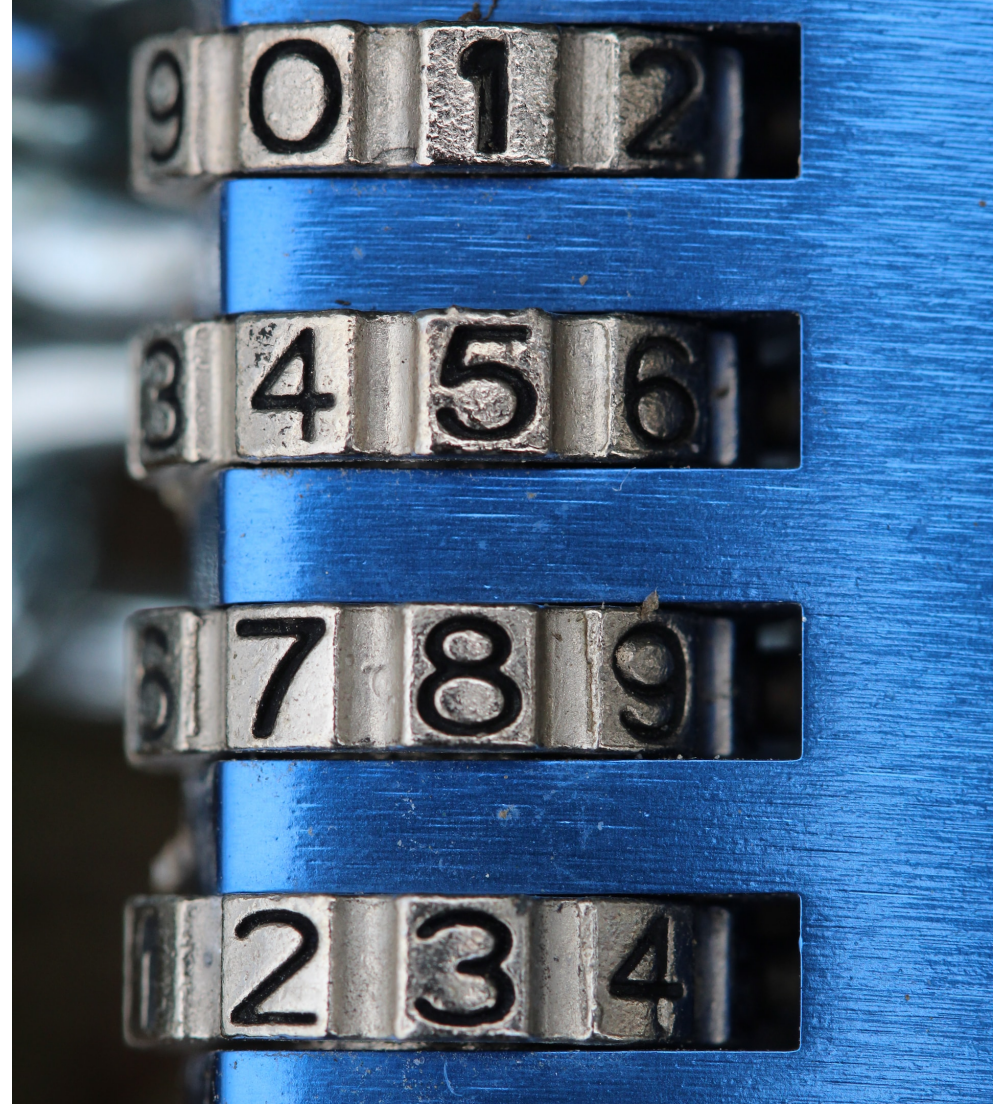
Des données ouvertes pour qui ?

- Chercheurs
- Entreprises
- Citoyens



Freins à l'ouverture

- Pas d'incitation à libérer ces données
- Les traces d'apprentissage sont des données identifiantes



FUN est un Groupement d'Intérêt Public

Notre crédo :

Argent public, code public ...
données publiques.

- Argent public ✓
- Code public * ✓
- Données publiques 🖥️

* <https://github.com/openfun>



Vers une ouverture de nos données d'apprentissage

Défis

1. Quel standard adopter pour ces données ?
2. Comment les anonymiser ?



Standard adopté : profils xAPI

- Un vocabulaire contrôlé pour des statements xAPI
- Fusion de datasets de sources différentes
- github.com/Gaia-X-DaSES

```
scripts: {  
  "start": "electron .",  
  "dev": "rollup -c -w",  
  "build": "rollup -c"  
},  
"keywords": [],  
"author": "",  
"license": "ISC",  
"devDependencies": {  
  "electron": "8.2.1",  
  "electron-reload": "1.5.0",  
  "concurrently": "5.1.0",  
  "@rollup/plugin-commonjs": "11.0.0",  
  "@rollup/plugin-node-resolve": "7.0.0",  
  "rollup": "1.20.0",  
  "rollup-plugin-livereload": "1.0.0",  
  "rollup-plugin-svelte": "5.0.3",  
  "rollup-plugin-terser": "5.1.2",  
  "svelte": "3.21.0"  
}
```

Anonymisation

En 2022, FUN est lauréat du Bac à Sable Education de la CNIL.

Premières étapes :

- Minimisation des traces d'apprentissage
- Hachage identifiant utilisateur
- Obfuscation de la clé de cours / ressource pédagogique



Anonymisation

Deuxième étape (avant publication) :

- Analyse statistique des traces
- Evaluer la génération de données synthétiques
- Participation à des challenges de ré-identification



Des questions, des remarques ?

Parlons-en.