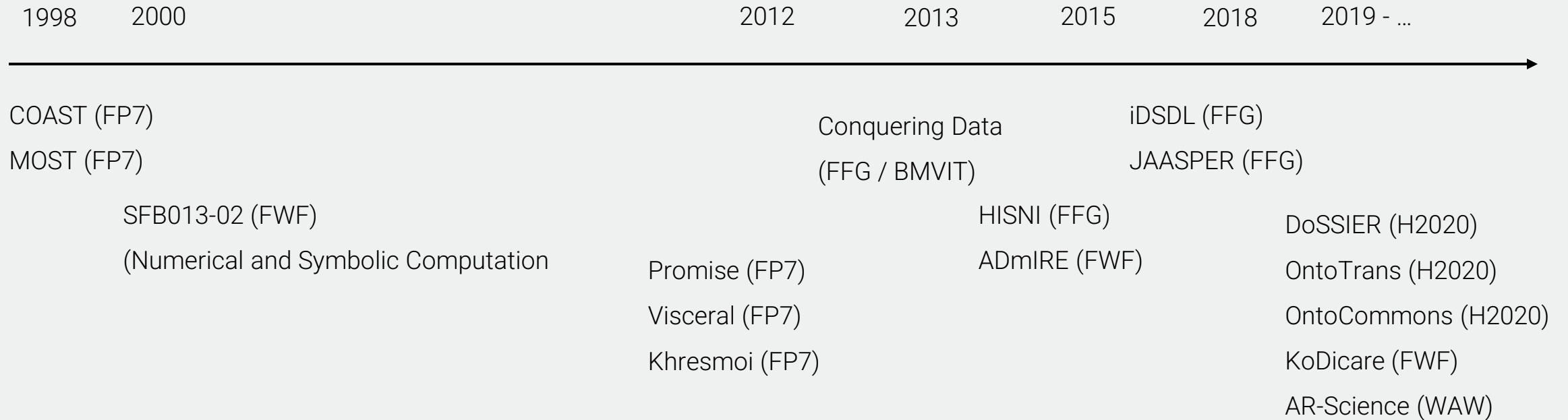# So many projects so little time

## … for DMPs …

## What do researchers need?

Florina Piroi

# Who am I

- Computer Scientist (since 1996)

- MSc In Parallel and Distributed Computing (1998)

- PhD in Symbolic Computation (2004) (Tools for Mathematical Knowledge Management)

- Researcher (since 2000)

  - Data Science

  - Domain Specific Information Retrieval

  - Machine Learning

  - Natural Language Processing

  - …

- Joined the Research Data Management Team at TU Wien couple of weeks ago.
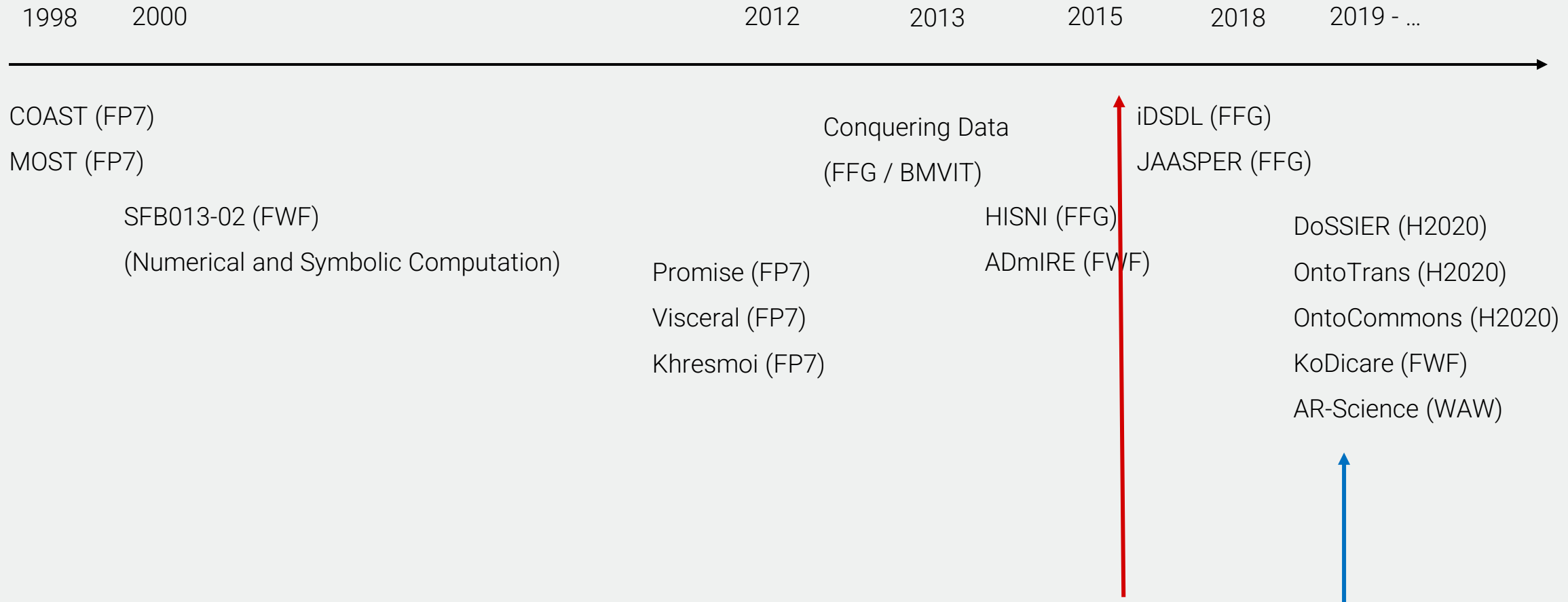
# Projects – a timeline

1998     2000                              2012         2013         2015       2018      2019 - …

COAST (FP7)

MOST (FP7)

Conquering Data

(FFG / BMVIT)

iDSDL (FFG)

JAASPER (FFG)

SFB013-02 (FWF)

(Numerical and Symbolic Computation

HISNI (FFG)

ADmIRE (FWF)

DoSSIER (H2020)

Promise (FP7)

Visceral (FP7)

Khresmoi (FP7)

OntoTrans (H2020)

OntoCommons (H2020)

KoDicare (FWF)

AR-Science (WAW)

# Data sets

- Medical content

- IP (Patent) documents

- Scientific Articles

- (social network) Logs

- Mathematical content (formulae, proofs, MathML formatted)

- Images (flow-charts, diagrams, chemical formulae, technical drawings

- Industry data (material science, interdisciplinary projects)

- Ontologies & Knowledge Graphs

- …

# Projects – a timeline

1998    2000                                    2012        2013        2015        2018    2019 - …

COAST (FP7)

MOST (FP7)

SFB013-02 (FWF)

(Numerical and Symbolic Computation)

Conquering Data

(FFG / BMVIT)

iDSDL (FFG)

JAASPER (FFG)

HISNI (FFG)

ADmIRE (FWF)

DoSSIER (H2020)

OntoTrans (H2020)

Promise (FP7)

Visceral (FP7)

Khresmoi (FP7)

OntoCommons (H2020)

KoDicare (FWF)

AR-Science (WAW)

# DMPs – Selected examples

iDSDL (FFG)

OntoTrans (H2020)

DoSSIER (H2020)

# iDSDL (FFG)

- Innovationslehrgang Data Science and Deep Learning

- Transfer of SoTA knowhow on AI, DL, DS to industries in Austria

- 20+ Industrial partners

- 5 teaching modules

- 1 Transfer project per industrial partner

  - Had to also create a DMP

- Expected 20+ DMPs – got lots lots fewer.

# iDSDL (FFG)

- Innovationslehrgang Data
- Transfer of SoTA knowho
- 20+ Industrial partners
- 5 teaching modules
- 1 Transfer project per indu
  - Had to also create a

## 5. Data Management Plan

### Data Summary

We use annotated image datasets to train and evaluat
segmentation and retrieval). The data is used to develop
developed during the course of the project. We collect fr
scale product catalogue (20+ million products, 2+TB of
such as OpenImages or ImageNet, or by web scraping of
are specific to the fashion retail industry, and includes i
clothing, beauty and accessories products.

### F.A.I.R. Data

**Findable**

We store datasets as close to their original format as po
same datasets. Our own annotations are mostly kept sep
developed in-house, transformations applied to the data
applied online during training, evaluation, etc.

This way the original data can be used without un
hundreds-of gigabytes sized datasets. Tiered memory and

**Accessible**

Most datasets created thus far by us have been proprietary for internal use only and could only be shared in rare exceptions, most of the rights to the images have been granted to use in a limited, non-transferable license agreement through contracts with online retailers. Other data that was used that might be permissively licensed, for instance under creative commons is already publicly available and in publications we do cite the corresponding works, where to acquire the datasets and/or where to acquire the necessary commercial licenses.

**Interoperable**

In our case that is mostly applicable to image annotations (labels). For some datasets we use the tooling for annotation, training/evaluation and deployment that was developed in-house and that is using formats that are somewhat proprietary as well. In other cases annotations are stored in industry-standard formats, such as COCO object detection bounding boxes and RLE (run-length encoding) for segmentation datasets.

All datasets we use are stored in open data exchange formats and data serialization formats that are free, open and interoperable. Other than TSV/CSV, JSON and JSON-Lines, we also use MessagePack and LMDB.

**Re-Use**

*(not applicable, see Accessible above)*

**Allocation of resources**

Most of the principles of F.A.I.R have already been fulfilled during the whole research and development phases since they already align with many industry best practices. Other aspects of open data do not apply in our case for reasons described under Accessible above.

Since we are a small start up, at the moment our data storage sizes are not significantly large, in the orders of 5-15 TB and only relatively slowly increasing. The cost of long term storage is therefore mostly negligible for us at this point, cold backup storage in the cloud is very cheap (only incurring

# iDSDL (FFG)

- Innovationslehrgang Data Science and De...
- Transfer of SoTA knowhow on AI, DL, DS...
- 20+ Industrial partners
- 5 teaching modules
- 1 Transfer project per industrial partner
    - Had to also create a DMP

---

## IDSDL - Data Management Plan

In order to calculate the expected revenue per location based on their environment we directly accessed th███████ datawarehouse. No personal data was used neither used or required.

Since there already was a project calculating the location potential based on a rule-based system, we re-used the data that was the basis for this calculation as well:

- Location data (historical)
    - Sales/Ticket data of relevan███████roducts████████████████████ the target dimension y
    - Branch type: 10 categories
    - Opening hours: start and end times per weekday, seasonal or not
    - Customer frequency: The average customer count during regular (non-seasonal) opening hours
    - ZIP
    - Contract partner: One contract partner/entity can own/manage multiple locations
    - Municipality-data (partly based on GEO-GIS data)
        - Population (/w and /wo commuters)
        - Purchasing power
        - Locations per municipality

Additionally we used the values generated by the rule-based system as additional features (allthough the effects on the models should correlate with the given base-data):

- Potential data as calculated by the other project: Relative, scaled between 1 and 5 partly based on fixed borders/rules
    - Branch potential (derived from banch type)
    - Customer frequency potential (derived from customer frequency)
    - Density potential for municipality (derived from population and number of locations)
    - Opening hours potential (based on opening hours)
    - Purchasing power potential (based on purchasing power)
    - Overall calculated potential given a set of pre-defined weights per potential metric

Since the data was already used by a project in production we did not spend too much time exploring (i.e. no outlier analysis) nor verifying the data.

# iDSDL (FFG)

- Innovationslehrgang Data Sc
- Transfer of SoTA knowhow
- 20+ Industrial partners
- 5 teaching modules
- 1 Transfer project per indust
  - Had to also create a D

**Data Collection**

**What data will you collect or create?**

Word und PDF Dateien die Verträge und Lohnabrechnungen darstellen. Vers

**How will the data be collected or created?**

Die Daten werden direkt von der Kanzlei übergeben.
Die Struktur der Daten ist nicht von Relevanz. Diese
alle Ordner ein.
Die Daten werden weiterhin auch im Produktivbetrie
tägliche Verwendung der Kanzlei werden.
Die Qualität wird aus der Natur des Geschäfts unsere

als Kunde muss sich darum keiner weiteren Sorgen machen. Alle Mitarbeiter der Kanzlei Algorithmus trainiert wird und das Projekt in Umsetzung ist. Danach muss nur noch die A

Created using DMPonline. Last modified 30 March 2020

**Documentation and Metadata**

**What documentation and metadata will accompany the data?**

Es handelt sich bei den Daten um Worddateien. Diese werden von mehreren
gestellt. Metadaten liegen nur in Form von Word-Metadaten vor und werden

**Ethics and Legal Compliance**

**How will you manage any ethical issues?**

Die Daten liegen in sicheren Cloudlösungen bereit. Es handelt sich dabei um

# OntoTrans (H2020)

- Ontology Driven Open Translation Environment

- Material Modelling / Material Sciences


- Ontology creation for specific domains (steel, chemical, prepregs / curing ) such that the translation scheme is efficient, cross-domain, adaptable, etc.

- Ontotrans.eu

# OntoTrans (H2020)

- DMP created more systematic, now

- EU H2020 template

- Extensive discussions with industrial partners what exactly this plan means in terms of:

  - Access rights and security

  - Disclosure of industrial proprietary data

  - Concept clarification (especially FAIR)

  - …

# OntoTrans (H2020)

- DMP created m

- EU H2020 temp

- Extensive discu

  terms of:

  - Access ri

  - Disclosur

  - Concept

# OntoTrans (H2020)

- DMP created more sys
- EU H2020 template
- Extensive discussions
  terms of:
  - Access rights
  - Disclosure of indu
  - Concept clarifica

All collected into one Confidential deliverable (pdf file!)

# DoSSIER (H2020)

- Domain Specific Systems for Information Retrieval

- Dossier-project.eu

- MSCA ITN/ETN

  - 15 subprojects (15 PhDs)

  - Lots! of data!

- Instructions sent to students with tables to fill

- Top-down approach

# DoSSIER (H2020)

- Domain Specific Systems for Information Retrieval

- Dossier-project.eu

- MSCA ITN/ETN

  - 15 subprojects (15 PhDs)

  - Lots! of data!

- Instructions sent to students with tables to fill

- Top-down approach

*Table 1 List of Data sets in DoSSIER*

| No. | Data set Name / Description | Partner(s) | Produced / Reused | Volume | License | Personal data |
|---|---|---|---|---|---|---|
| 1 | CLEF-IP – patent documents | TUW, IHU ESR2, ESR4 | Reused | 14GB | CC-NC-SA 3.0 | No |
| 2 | MAREC / IREC – patent document collection | IHU, ESR2 | Reused | 621GB | CC-NC-SA 3.0 | No |
| 3 | AC_1 (working title) | USFD, ESR3 | Produced | ~10GB | TBD | Not stored |
| 4 | COLIEE - Competition on Legal Information Extraction/Entailment | TUW, ULEI ESR4, ESR6 | Reused | 2GB | Free for research | No |
| 5 | CaseLaw – legal case documents | TUW, ESR4 | Reused | | Free for research | No |
| 6 | TripClick – click log data set | TUW, ESR4 | Reused | 32GB | Free for research | No |
| 7 | TREC-COVID | TUW, ESR4 | Reused | | Free for research | No |
| 8 | MS Marco | TUW, ESR4 | Reused | | Free for research | No |
| 9 | ES_1 (working title) | SUG, ESR5 | Produced | ~1GB | TBD | Not stored |
| 10 | SciDocs – Scientific Documents | ULEI, ESR6 | Reused | | Free for research | No |
| 11 | CIR_1 (working title) | UMB, ESR8 | Produced | ~10GB | TBD | Not stored |
| 12 | TE_1 (working title) | USFD, ESR10 | Produced | ~10GB | TBD | Not stored |
| 13 | MCC_1 (working title) | SUG, ESR11 | Produced | ~100MB | TBD | Not stored |
| 14 | EMIS_1 (working title) | SUG, ESR12 | Produced | ~700GB | TBD | Not stored |
| 15 | TREC 2021 Clinical Trials | UMB, ESR14 | Reused | | Free for research | No |
| 16 | NFCorpus | UMB, ESR15 | Reused | 27MB | Free for research | No |
| 17 | Legal_data | UMB, ESR15 | Reused | 2GB | Free for research | No |

# DoSSIER (H2020)

*Table 1 List of Data sets in DoSSIER*

| No. | Data set Name / Description | Partner(s) | Produced / Reused | Volume | License | Personal data |
|---|---|---|---|---|---|---|
| 1 | CLEF-IP – patent documents | TUW, IHU ESR2, ESR4 | Reused | 14GB | CC-NC-SA 3.0 | No |
| 2 | MAREC / IREC – patent document collection | IHU, ESR2 | Reused | 621GB | CC-NC-SA 3.0 | No |
| 3 | AC_1 (working title) | USFD, ESR3 | Produced | ~10GB | TBD | Not stored |
| | COLIEE - Competition on Legal | | | | | |

- Domain Specific Systems for Information Retrieval

- Dossier-project.eu

- MSCA ITN/ETN
  - 15 subprojects (15 PhDs)
  - Lots! of data!

- Instructions sent to students with tables to fill

- Top-down approach


- Each data set described separately,

  in an additional table

| Data Set No. 1 | |
|---|---|
| **NAME or Identifier** | **CLEF-IP** |
| **DoSSIER Project/ESR** | P02 / ESR2, P04 / ESR4 |
| **Description** | A collection of more than 1.3M patent documents (~2.6 million files) derived from EPO (European Patent Office) sources and EuroPCT Applications (more than 400K documents) published by WIPO (World Intellectual Property Organization). The collection contains documents in English, French and German with at least 150,000 documents in each language, all published before 2001. |
| **Re-used Data** | Yes |
| **Standards and Metadata** | Dtd available |
| **File Format** | XML |
| **Size** | 14Gb |
| **Data Sharing** | Open |
| **Access Rights** | CC-NC-SA 3.0 |
| **Archiving and Preservation** | https://researchdata.tuwien.ac.at/records/khw86-rnf37 |
| **Ethics & Legal Compliance** | Not the case |
| **Person Identifiable Data** | Not contained |

| 17 | Legal_data | UMB, ESR15 | Reused | 2GB | research | No |

# Lessons Learned

Did you notice?

- All just pdf files – static
- No follow up on them (that I am aware of)
- Not retrievable for statistics
- Very different information (template dependent)
- "why" not clear enough in the community
- KISS – researchers don't want overhead related to data management (but no way out – how do we tell them that?)
- And: often researchers (in my domain) release some data, somewhere (e.g. Kaggle, hugginface), little overlap with DMP (tool-ed or pdf-ed)

Informatics DATA SCIENCE 24.11.2022

# What we'd like

- Simple guidelines (actually available!)

- Small overhead

- Early introduction to the whole ecosystem of data management

- Proof of benefit

- Institutional support (advisors, community, etc)