SWEDEN

Data Readiness Lab

Databeredskapsverkstaden

Referensgruppmöte (2022-06-08)

Background

October 2021 – October 2023 Funded by

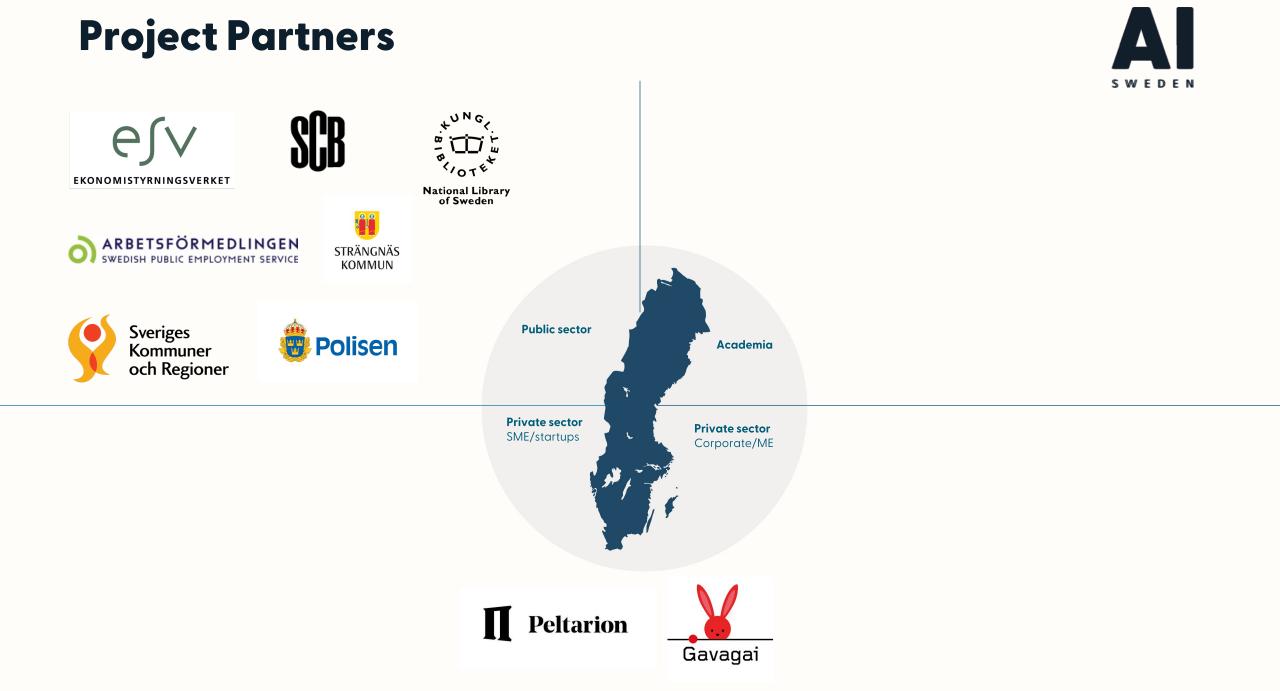
Data Readiness is

- an overlooked problem
- often poorly understood
- a huge bottleneck for the application of AI (Swedish Language Models)

Project focuses on

- the public sector
- text data













You want to apply a language model, but it turns out the data..

- .. is hard to get access to
- .. is not in the format you expected
- ... does not meet your expectations w.r.t. quantity
- .. does not meet your expectations w.r.t. quality (noise, missing values, ..)



- Availability
- Validity
- Utility





You want to apply a language model, but it turns out the data..

- .. annotation quality is insufficient (=> poor model quality)
- .. annotation takes too much time (= data is not enough)



Annotation

- Quality
- Efficiency
- Reusability





You want to apply a language model, but it turns out the data..

• ... contains personal information



Anonymization

Hands-on Application





You want to apply a language model, but it turns out the data..

• the later you find out, the worse it is!



Evaluation

- Data Quality Assessment
- Test Data





Evaluation

- Test Data
- Quality Assessment

Anonymization

Hands-on
Application

Data Readiness

- Availability
- Validity
- Utility

Annotation

- Quality
- Efficiency
- Reusability

Project Goal



Att i stor skala möjliggöra ökad databeredskap hos behovsägare genom att tillhandahålla verktyg, ramverk och resurser.

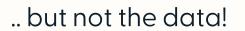




We share

- problems
- solutions

Data Readiness Lab











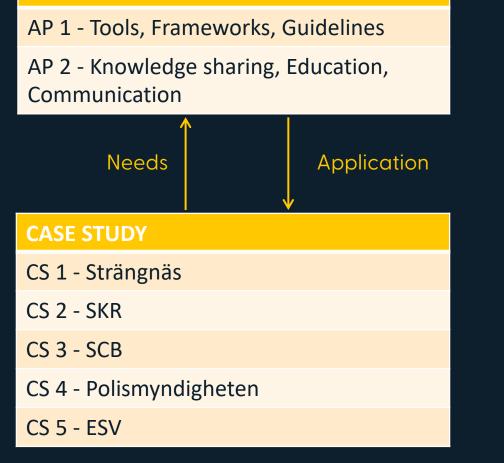
National Library of Sweden

TOTET

Work packages & Case Studies



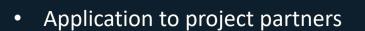
WORK PACKAGE



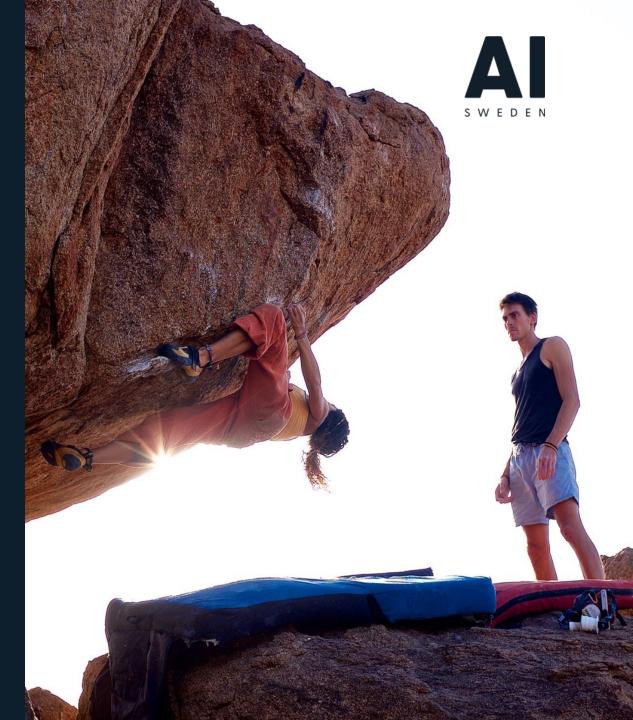
Data Readiness Annotation Anonymization Evaluation

First Step: Data Readiness Assessment

 Method: "We need to Talk About Data: The Importance of Data Readiness in Natural Language Processing" (Olsson & Sahlgren - 2021)





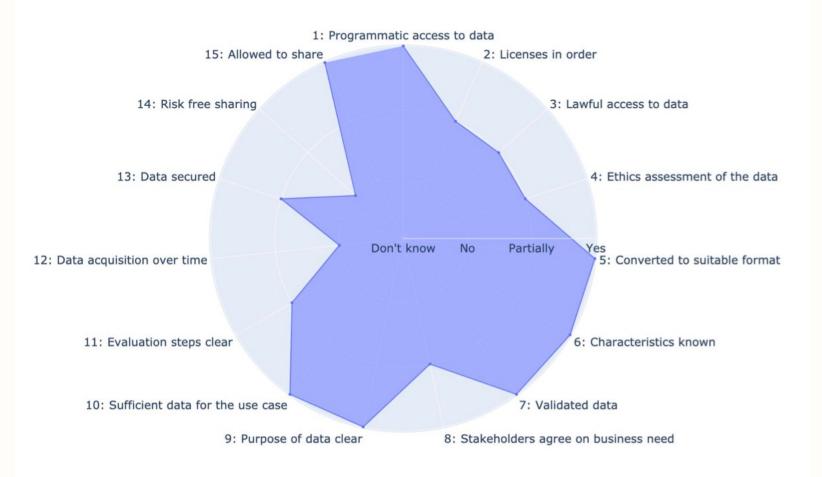




Go	bal	Band	Level	Meaning	The state of data
		Α	A-1	Utility	Ready for analysis with respect to given business objectives, hypotheses, questions, or context.
		В	B-1	Validity	Ready to define the candidate questions and hypotheses. Includes exploratory analysis, data characterization, entity disambiguation, and de-duplication, etc.
		С	C-1	Accessibility	Ready to be loaded in analysis software. Includes programmatic access, format conversions, and legal aspects, etc.
			•••		
Sta	art		C-n	Hearsay data	"I'm sure AI can solve this problem. We have lots of data!"
Lawrence (2017) <u>"Data Readiness Levels"</u> .					Gavagai

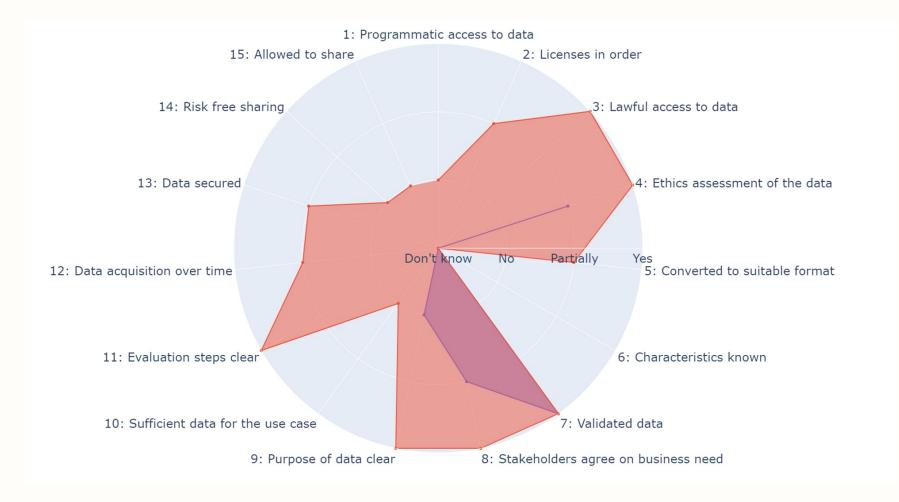






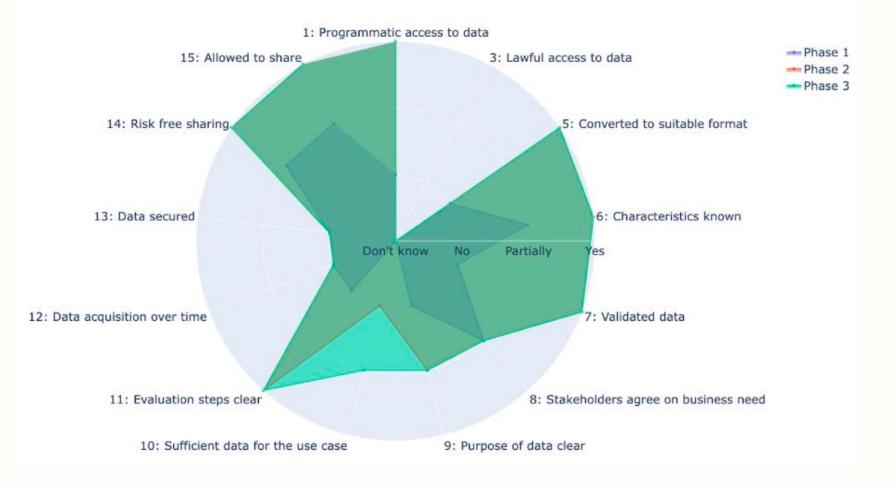














Peltarion





First Step: Data Readiness Assessment

- Method: "We need to Talk About Data: The Importance of Data Readiness in Natural Language Processing" (Olsson & Sahlgren - 2021)
- Application to project partners

Next Steps:

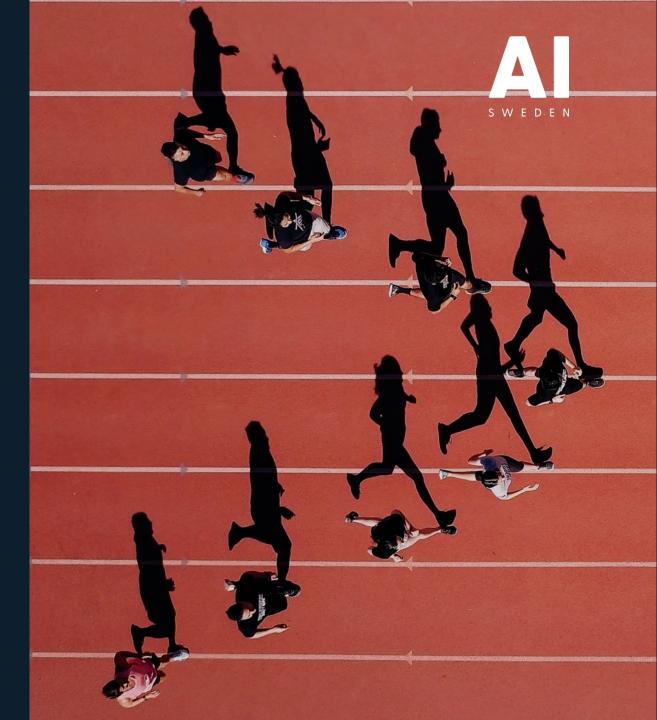
- Identification of major needs
- Creation of guidelines and tools for improved data readiness

First Step: Audit of Tools & Methods

Project Partners: Use Cases & Experiences

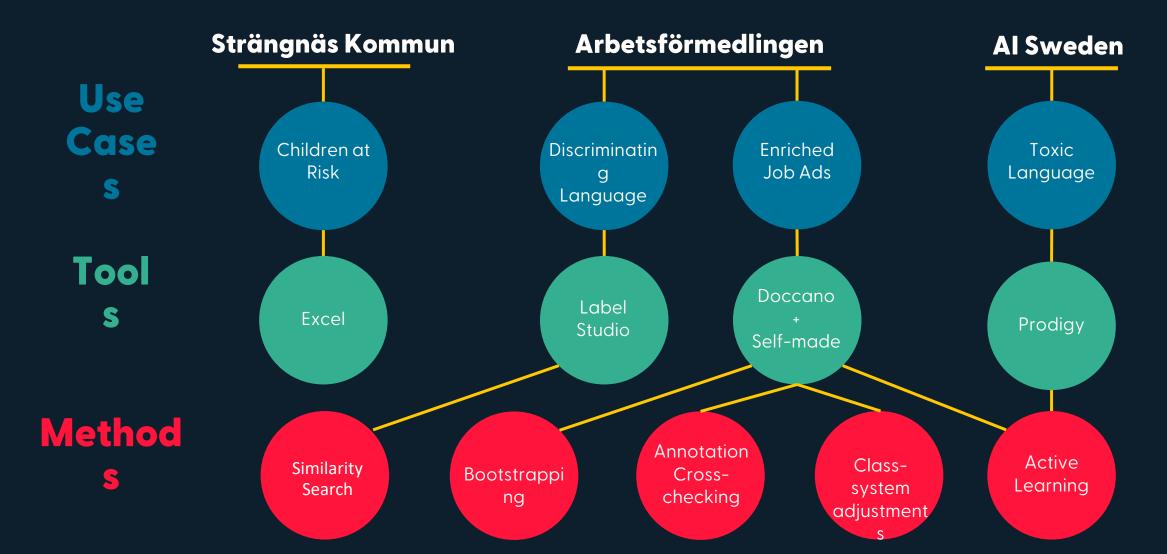
Feature Comparison of Text Annotation Tools (Olsson, https://nlp-data-readiness.readthedocs.io)

Job Ad Annotation at Arbetsförmedlingen (Stollenwerk, Fastlund, Nyqvist, Öhman – 2021)





Project Partners: Use Cases & Experiences





Feature Comparison of Text Annotation Tools (Olsson, https://nlp-data-readiness.readthedocs.io)

Tool	Prodigy	Label studio	Doccano
Website	https://prodi.gy/	https://labelstud.io/	https://doc
Tagline	"Radically efficient machine teaching. An annotation tool powered by active learning."	"Open Source Data Labeling Tool — Simplicity built-in, no overcomplicated UIs — Supports different datatypes — Visually configurable from top to the bottom"	"Text anno Just create data and s
SaaS	No	No	No
Self hosted	Yes, local server with web interface for annotation.	Yes	Yes
Active Learning	Yes, NER, text classification, dependecy parsing, pos tagging	No, but it can probably be implemented.	No, but it o
Multiple annotators for same data	Yes	No, but maybe.	Yes



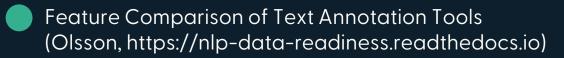
Job Ad Annotation at Arbetsförmedlingen (Stollenwerk, Fastlund, Nyqvist, Öhman – 2021)

METHOD	PURPOSE
Bootstrapping	Acceleration
Active Learning	Acceleration
Annotation Cross-checking	Quality Assurance
Class system Adjustments	Quality Assurance





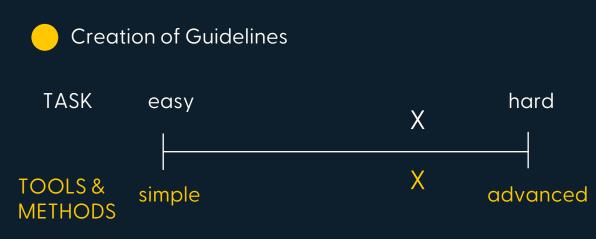
Project Partners: Use Cases & Experiences



Job Ad Annotation at Arbetsförmedlingen (Stollenwerk, Fastlund, Nyqvist, Öhman – 2021)

Next Steps:







Contact:

felix.stollenwerk@ai.se