AI SWEDEN

# Building GPT-SW3: The first LLM for the nordic languages.
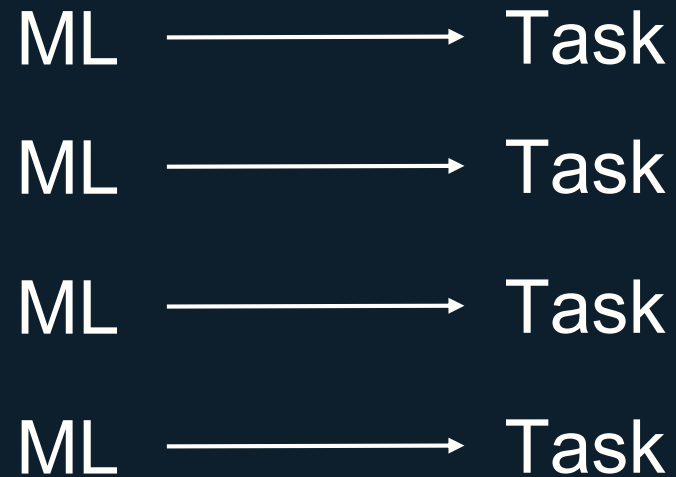
Language models learn language by reading text

Artificial neural networks (Transformers)
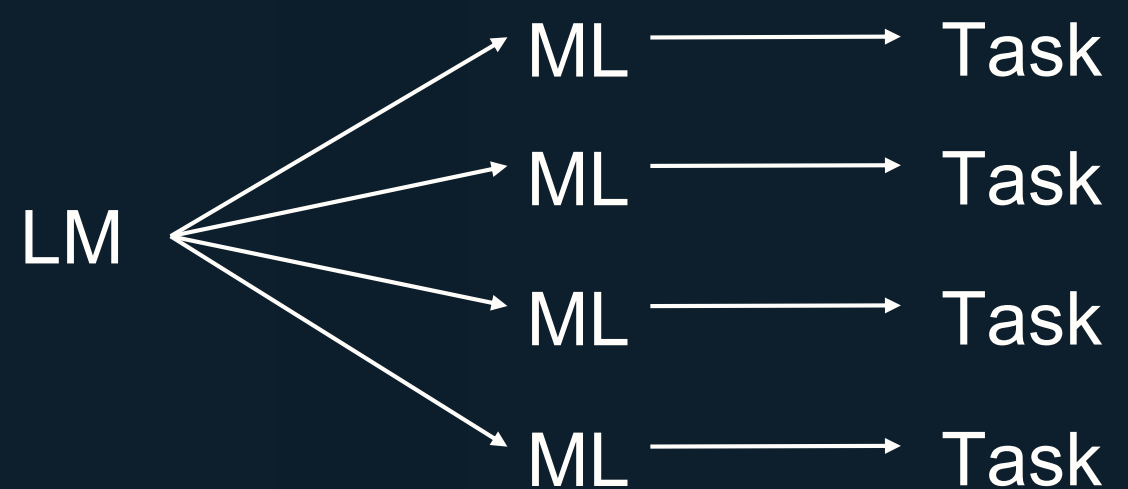
Can be used to solve (all) language processing tasks
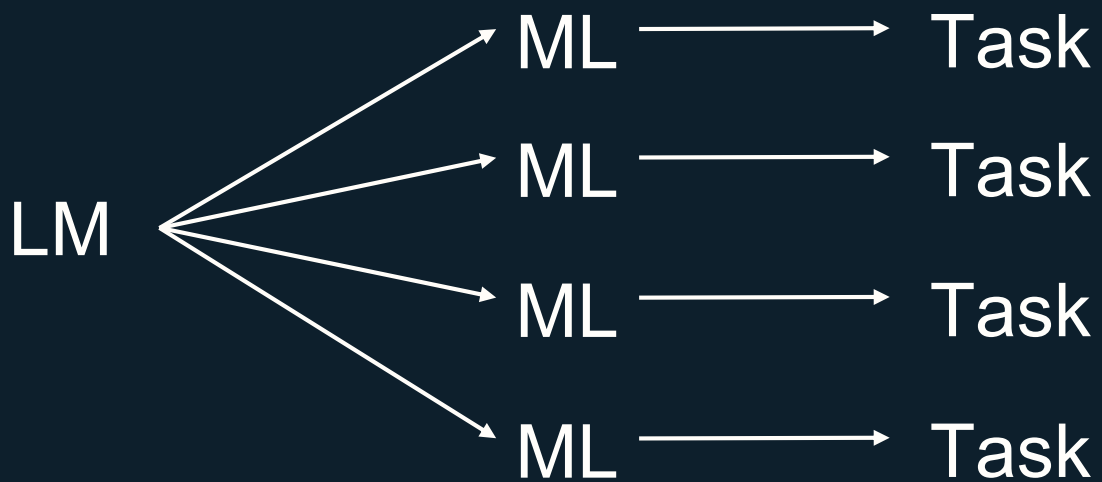
# Language models for Swedish authorities

Provide the prerequisites for authorities to use LMs
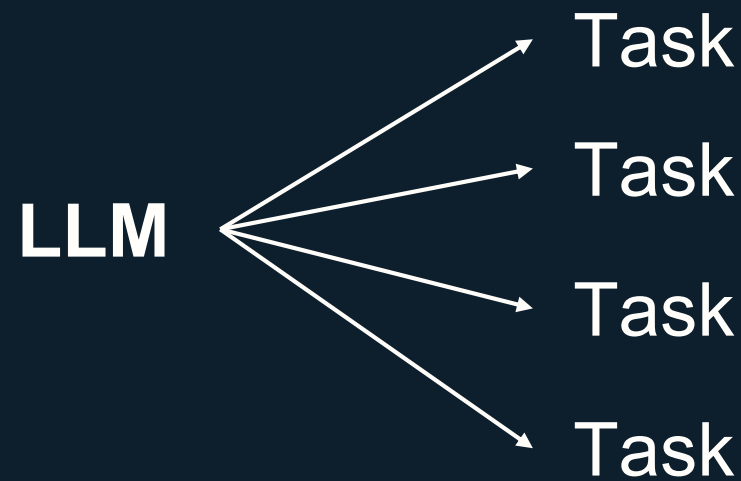
Code, data, models and applications

Most Swedish LMs have some connection to the project

All authorities in the project have built solutions based on LMs (Skatteverket, Arbetsförmedlingen, Tillväxtverket)
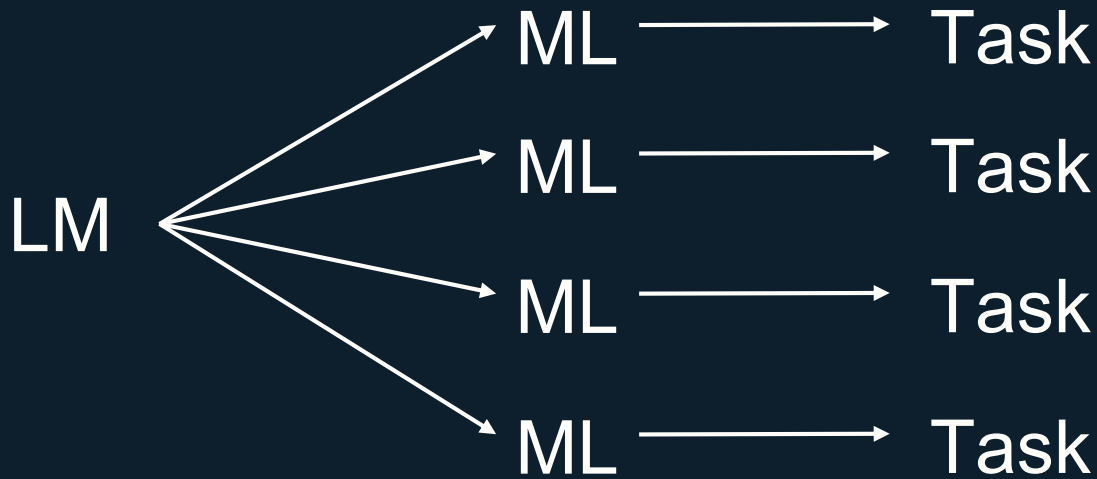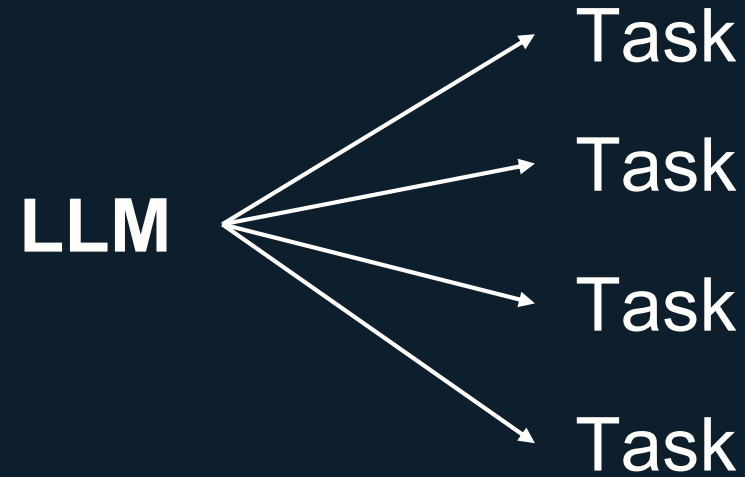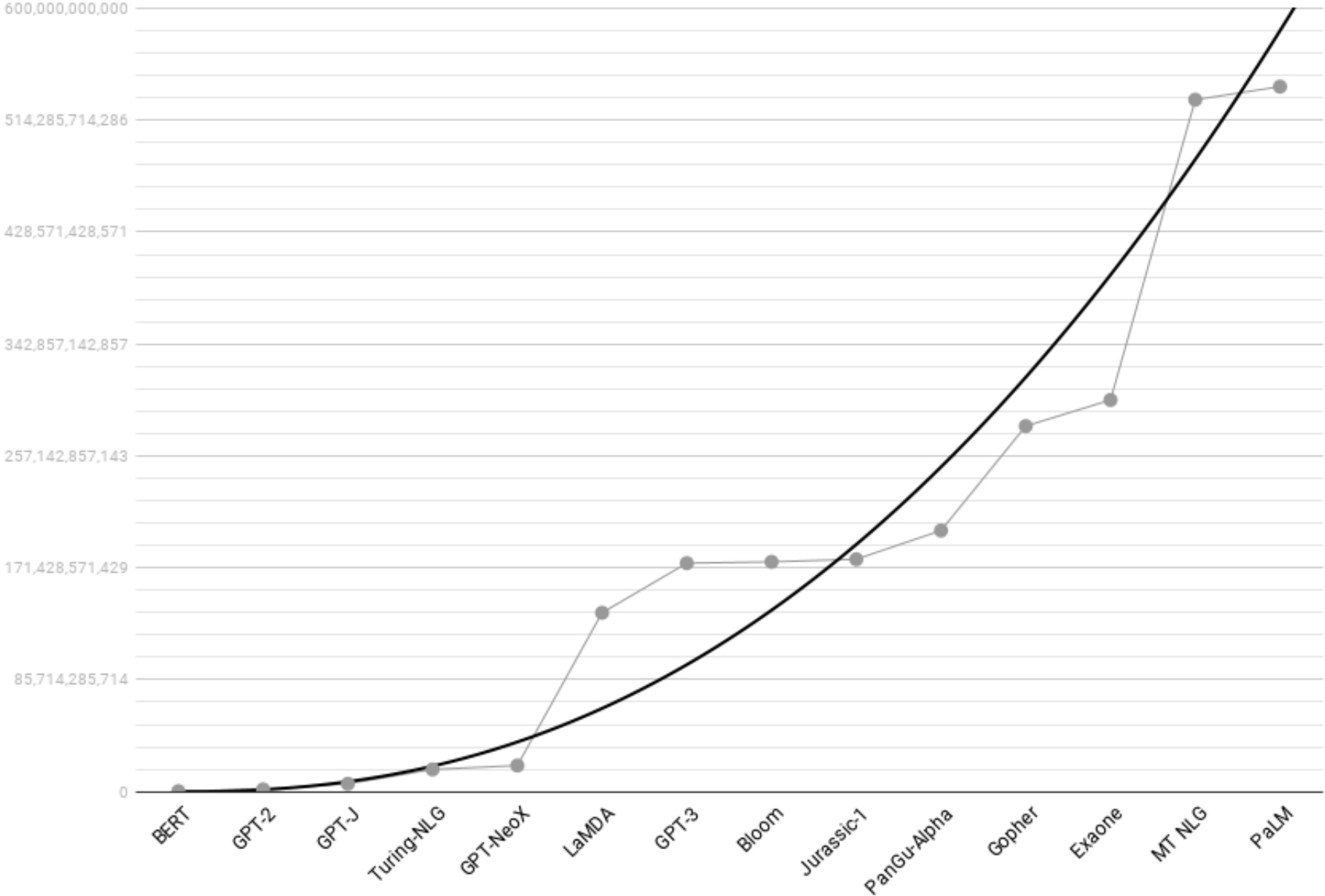
# Transformer model parameters

Bert large (340m): ~1000 GPUh
GPT-SW3 (6.7b): ~27000 GPUh
GPT-SW3 (40b): ~280000 GPUh

Bert large (340m): ~1000 GPUh ~ 5 DGX-days
GPT-SW3 (6.7b): ~27000 GPUh ~ 140 DGX-days
GPT-SW3 (40b): ~280000 GPUh ~ 4 DGX-years

source: https://developer.nvidia.com/blog/training-bert-with-gpus/
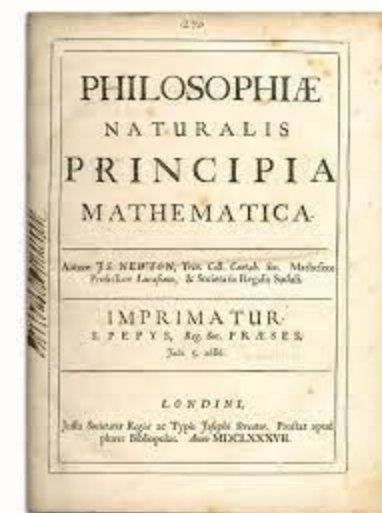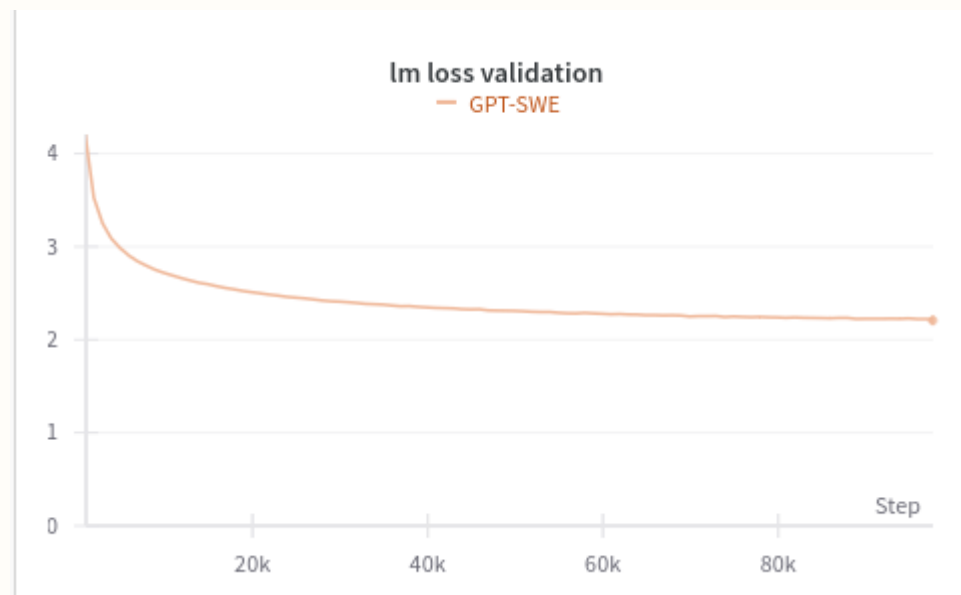
**NVIDIA SuperPOD:**
**60 NVIDIA DGX-A100 compute nodes**

**Each DGX-A100:**
**8 NVIDIA A100 Tensor Core GPUs**

**Each A100 GPU:**
**40 GB on-board HBM2 VRAM**

**Nvidia Mellanox Infiniband networking**

**Nvidia NeMo Megatron software stack**
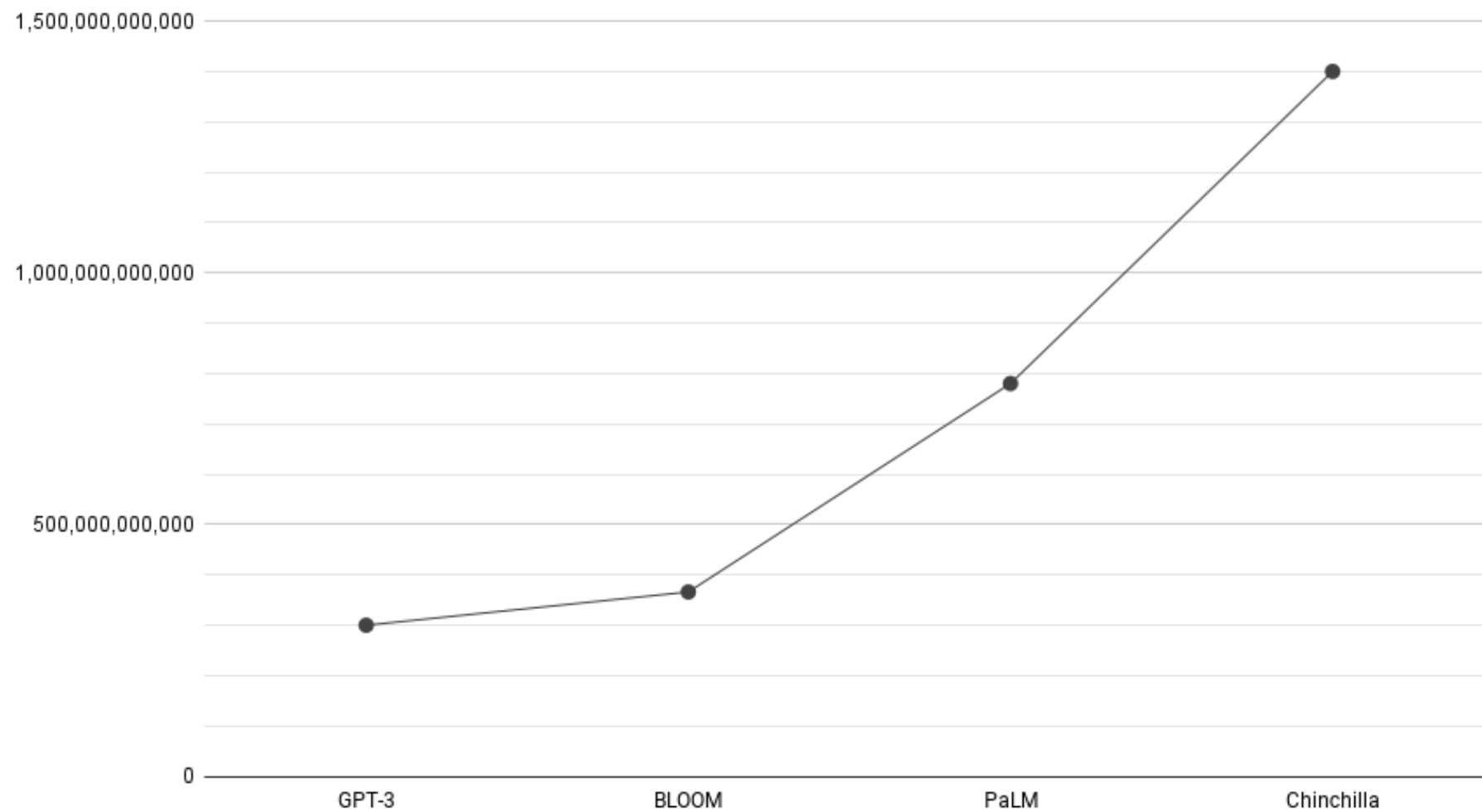
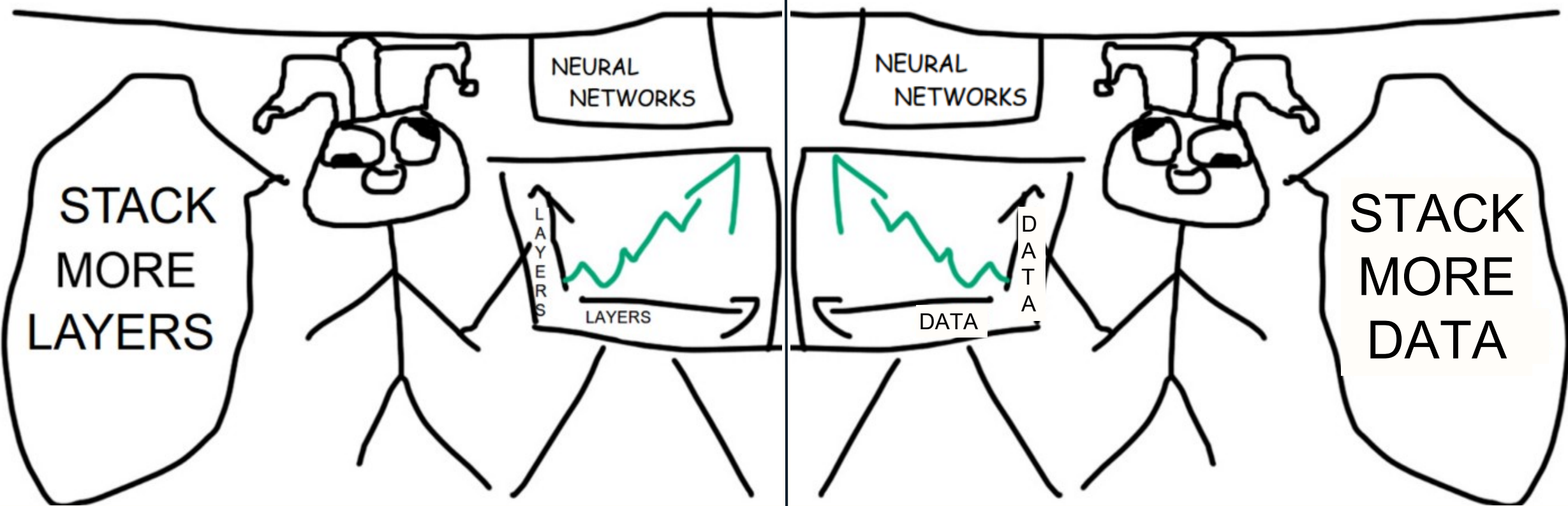## lm loss validation

— GPT-SWE

# GPT-SWE

- A 1.3 B parameter model trained using Megatron-LM
- Trained on found data: Oscar, MC4, Swedish parliamentary debates, et.c.
- Familiarity with HPC

Kaplan et al. (2020) Scaling laws for neural language models

Hofman et al. (2022) Training compute-optimal large language models

# GPT-SW3: a foundational resource for Nordic NLP

magnus.sahlgren@ai.se

# Why?

Build competence

Digital and linguistic sovereignty

Democratize access to LLMs
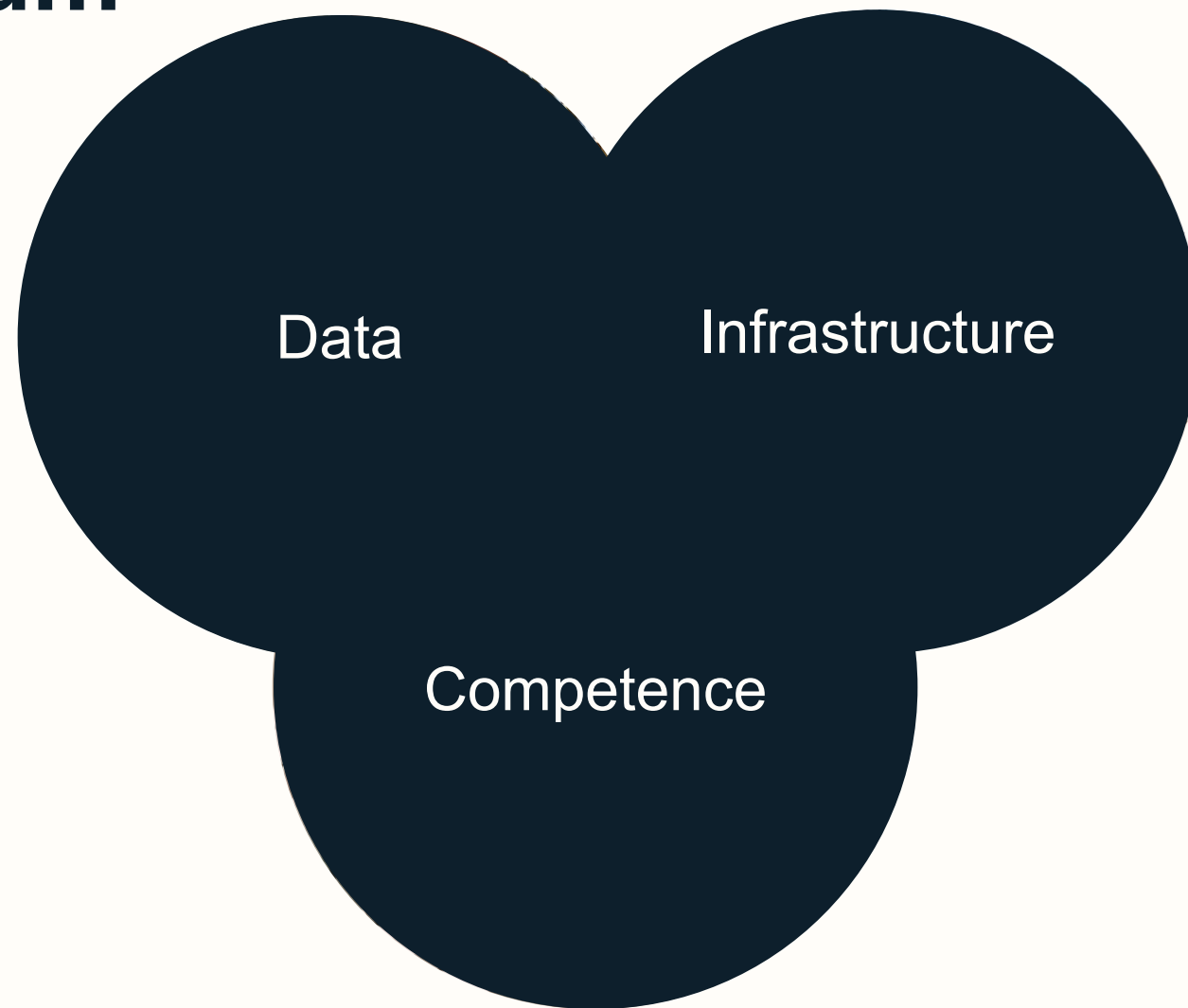
# How?

Collaboration

Transparency

# GPT-SW3 consortium

AI Nordics community on Discord

Data platform and deployment partners:

Data    Infrastructure

Competence

**AI**
SWEDEN

Swedish is a small language, but with close friends

Pool resources from typologically similar languages

The Nordic Pile: Swedish, Norwegian, Danish, Icelandic, English

Based on existing corpora and open sources
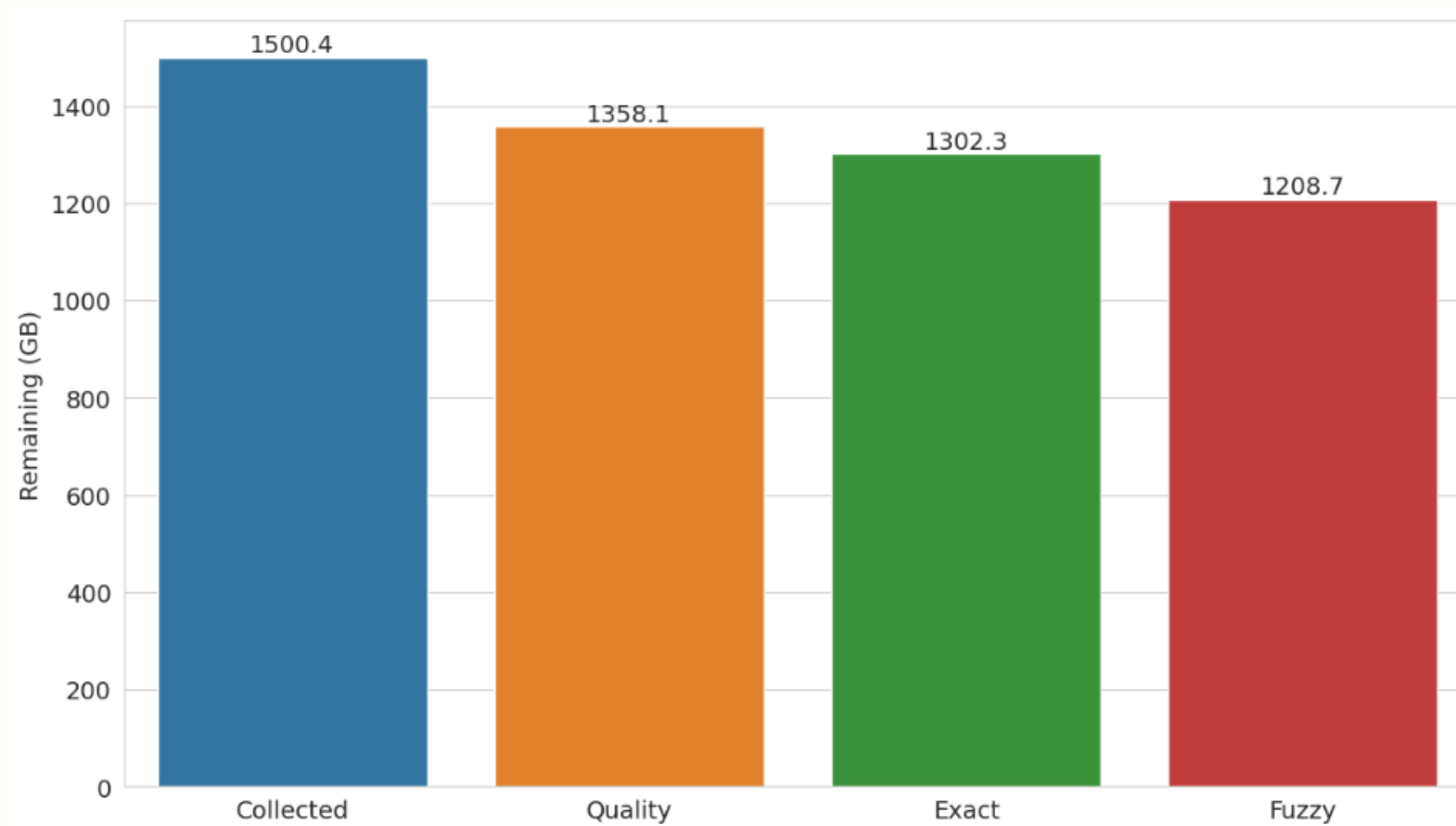
- **Normalization (whitespace, character encoding)**
- **Quality filtering**
- **Exact deduplication**
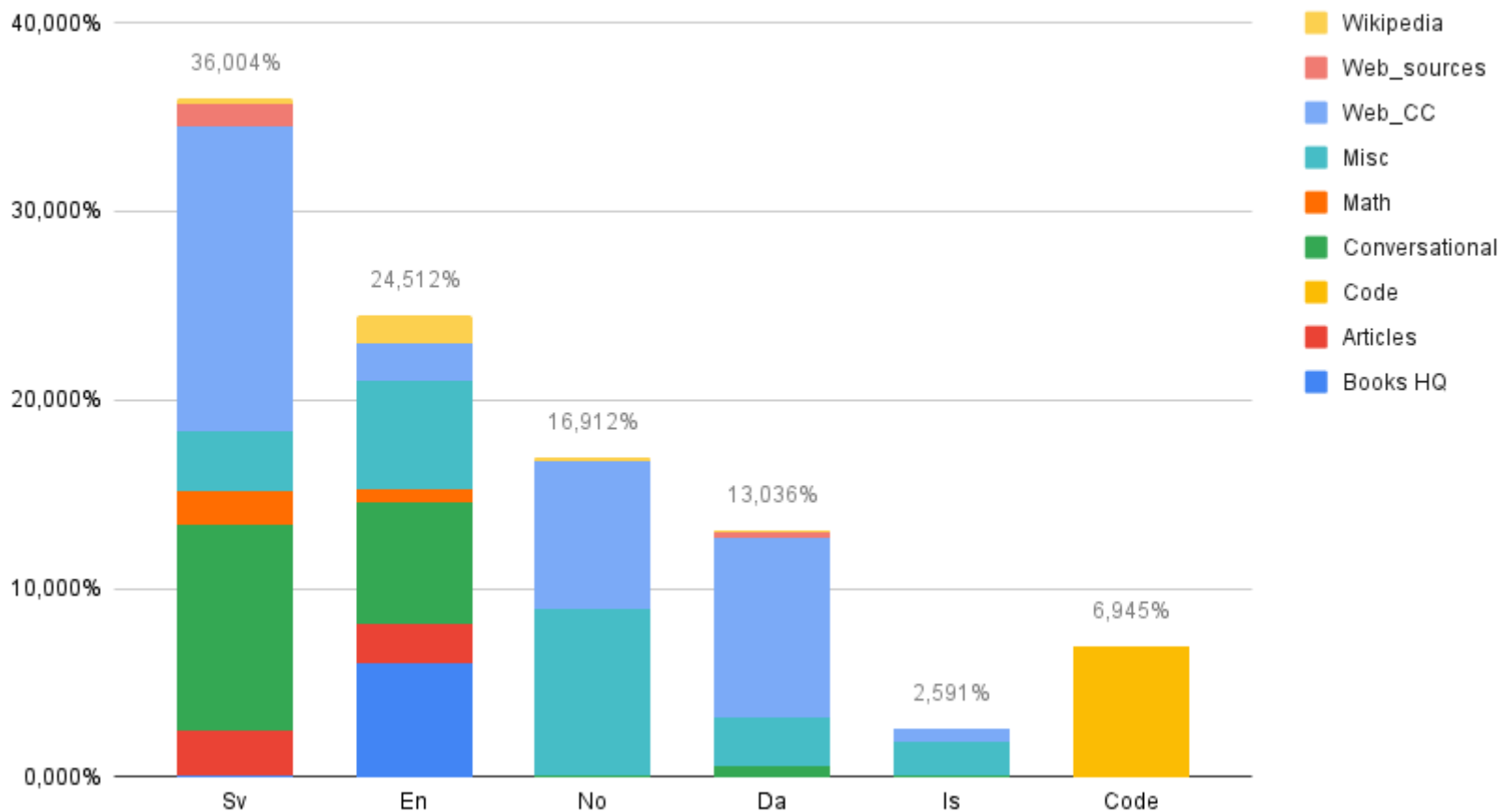- **Fuzzy deduplication (intra-language LSH)**
- **Weighting of data sources**

**Tokenizer: SentencePiece**
**Vocabulary size: 64k**

Sv, En, No, Da, Is …

Legend:
- Wikipedia
- Web_sources
- Web_CC
- Misc
- Math
- Conversational
- Code
- Articles
- Books HQ

| Category | Value |
|---|---|
| Sv | 36,004% |
| En | 24,512% |
| No | 16,912% |
| Da | 13,036% |
| Is | 2,591% |
| Code | 6,945% |

**January-March 2023:
6 pre-trained models**

126M – 40B parameter models

**April 2023:**

**4 instruction finetuned models added**

126M – 20B parameter models

**AI**
S W E D E N

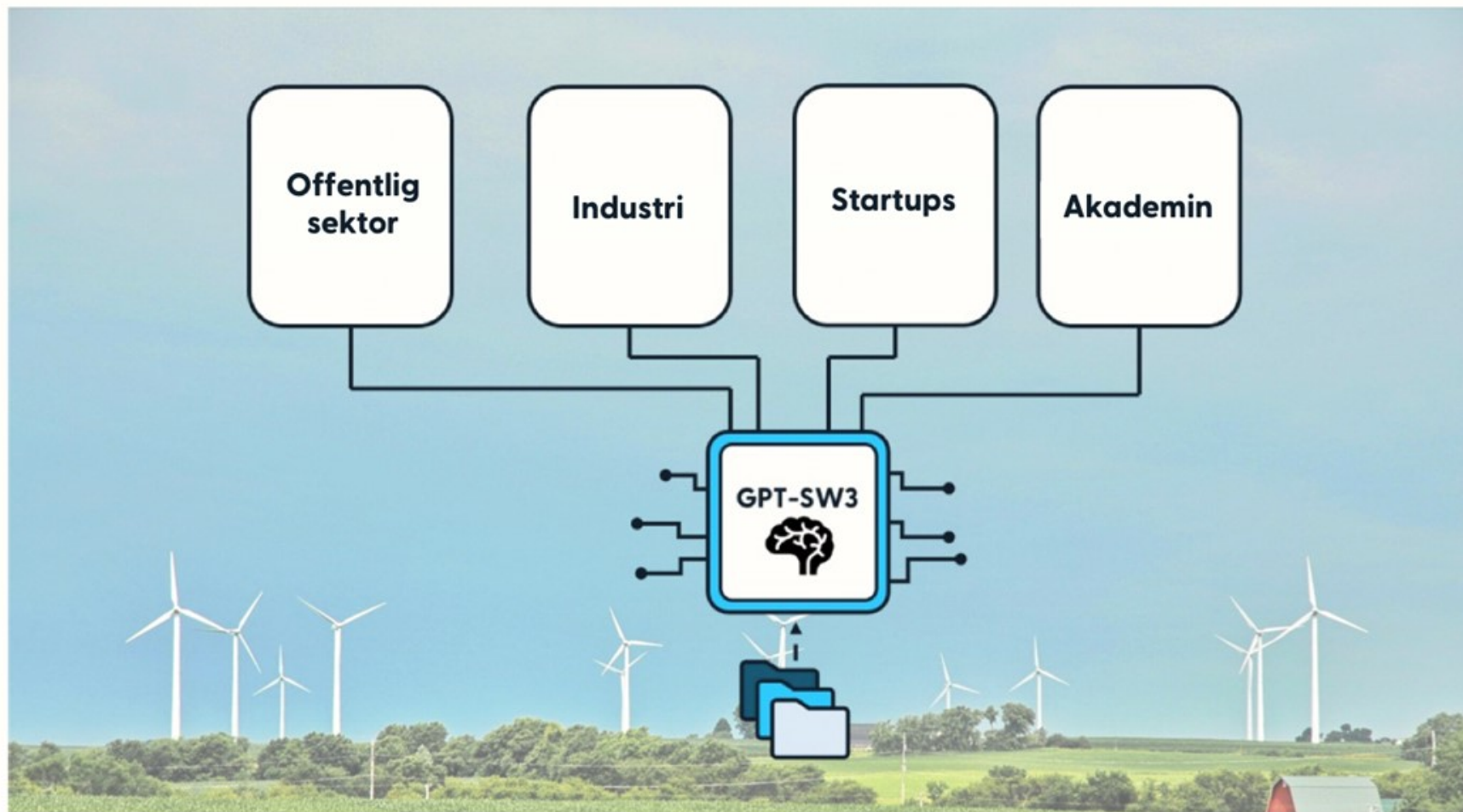| DESCRIPTION | EXAMPLE |
|---|---|
| **HEADER** | Tema: sol, moln, solstrålar, solnedgång |
| **INTRO** | Detta är en poetisk konversation mellan filosofen Josie och roboten Klara. De går på djupet om vad mänsklighet innebär och ger alltid varandra eftertänksamma, respektfulla och kloka svar.... |
| **THEMED QUESTIONS & ANSWERS** | Filosofen Josie: Hur får man energi?<br><br>Robotexperten Klara: Jag fylls av energi och vänlighet när jag tar del av Solen och hans näring. Men jag tror att nyckeln är att alltid sträva efter att vara den bästa versionen av sig själv, och att visa det för andra. Det är det enda sättet att skapa en bättre värld. |
| **USER QUESTION** | Filosofen Josie: [INSERT QUESTION HERE]<br><br>Robotexperten Klara: |

Initial prompt

# GPT-SW3: en svensk basmodell för texthantering, finansierad av Vinnova 2022-2024

Validate the use of GPT-SW3 for solving NLP tasks

Models, API, and applications

Private sector (small and large), public sector, academia

(Limited) serving of models

Instruction tuning of models

Validation of models "in anger"

**Feedback from participant
in pre-release**

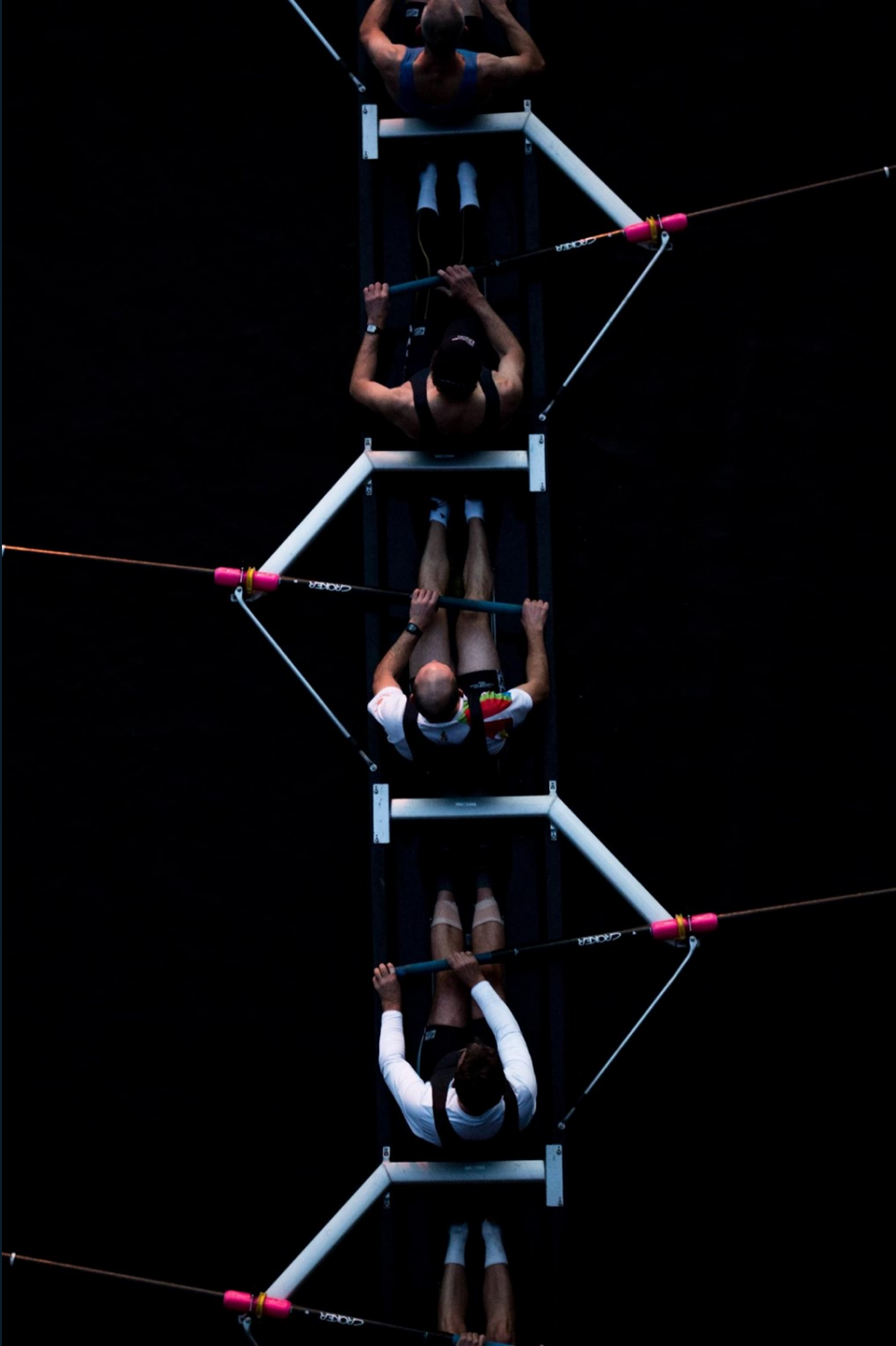" The democratisation of LLM is a true before-after moment. Thank you."

The way forward:

**Collaboration**

**Transparency**

**Investment**

# 🎉 GPT-SW3 Model release 🎉

Restricted pre-release January 2023
Open release November 16 2023 🎉

`ai.se/en/gpt-sw3`

magnus.sahlgren@ai.se