



TARTU ÜLIKOOL

The extent of legal control over language data: the case of language technologies

The presenters:

Aleksei Kelli

Penny Labropoulou





The authors:

Aleksei Kelli, Arvi Tavast, Krister Lindén, Kadri Vider, Ramūnas Birštonas, Penny Labropoulou, Irene Kull, Gaabriel Tavits, Age Värv



The aims of the research:

- To determine the extent of impact on the development of language technologies (LTs);
- To define a favourable regulatory framework for the development of LTs.

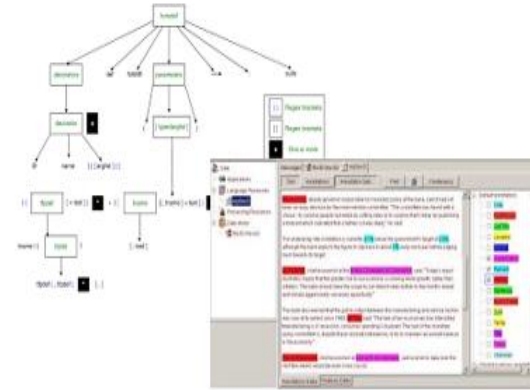
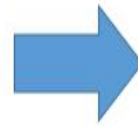
From Language Data to LTs:



raw data



datasets



annotated data



end-user applications

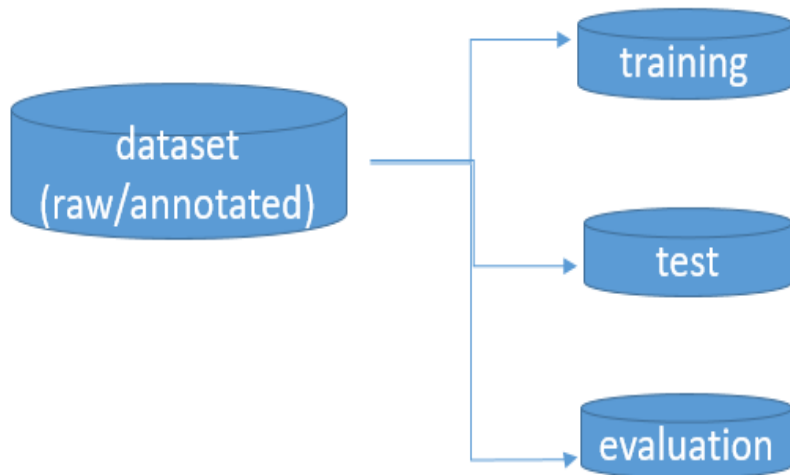


models

| Reddit temporal n-gram corpus | word | frequency | year | month |
|-------------------------------|-----------------------------|-----------|------|-------|
| 1-gram | trump | 981 | 2015 | 01 |
| 2-gram | trump apprentice | 31 | 2015 | 01 |
| 3-gram | donald trump battle | 16 | 2015 | 01 |
| 4-gram | donald trump ignorant tweet | 8 | 2015 | 01 |
| 5-gram | take donald trump advice in | 2 | 2015 | 01 |

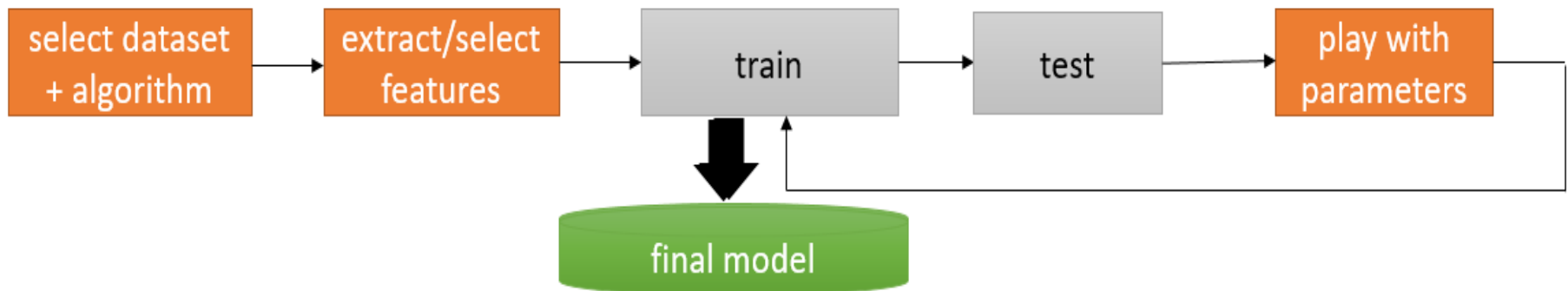
A diagram of a neural network with three layers of nodes. Below the table is a screenshot of a data visualization tool showing a list of items with various attributes and values.

Building models process:



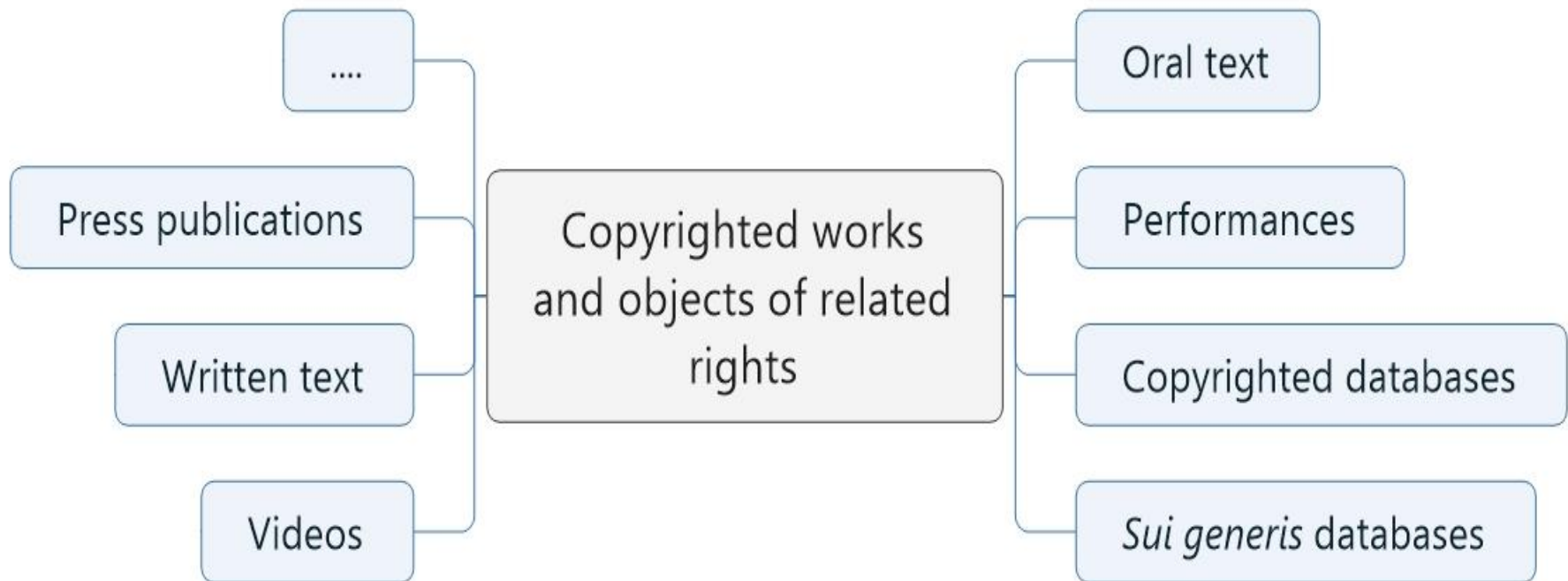
human-driven

auto processing

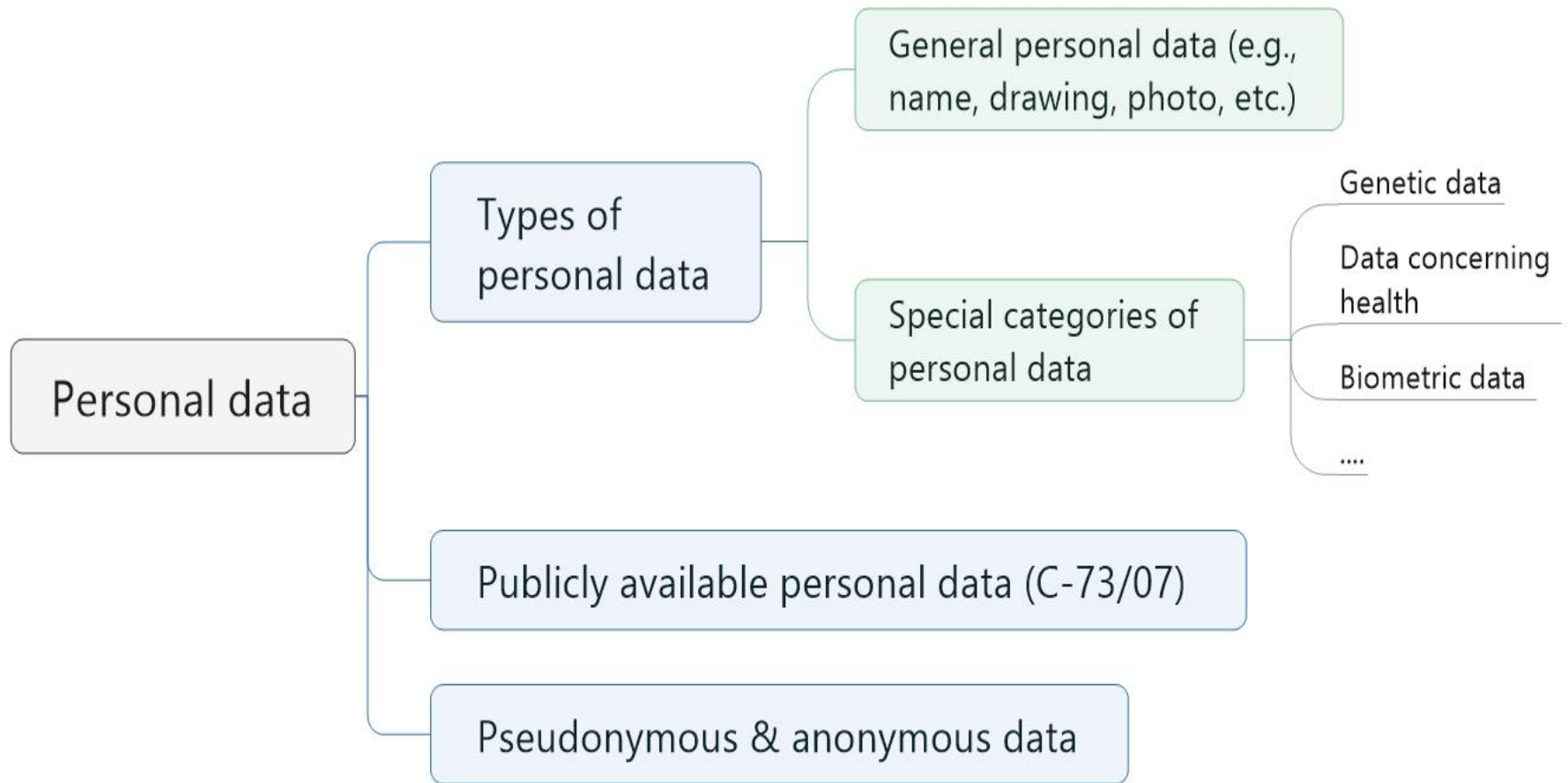




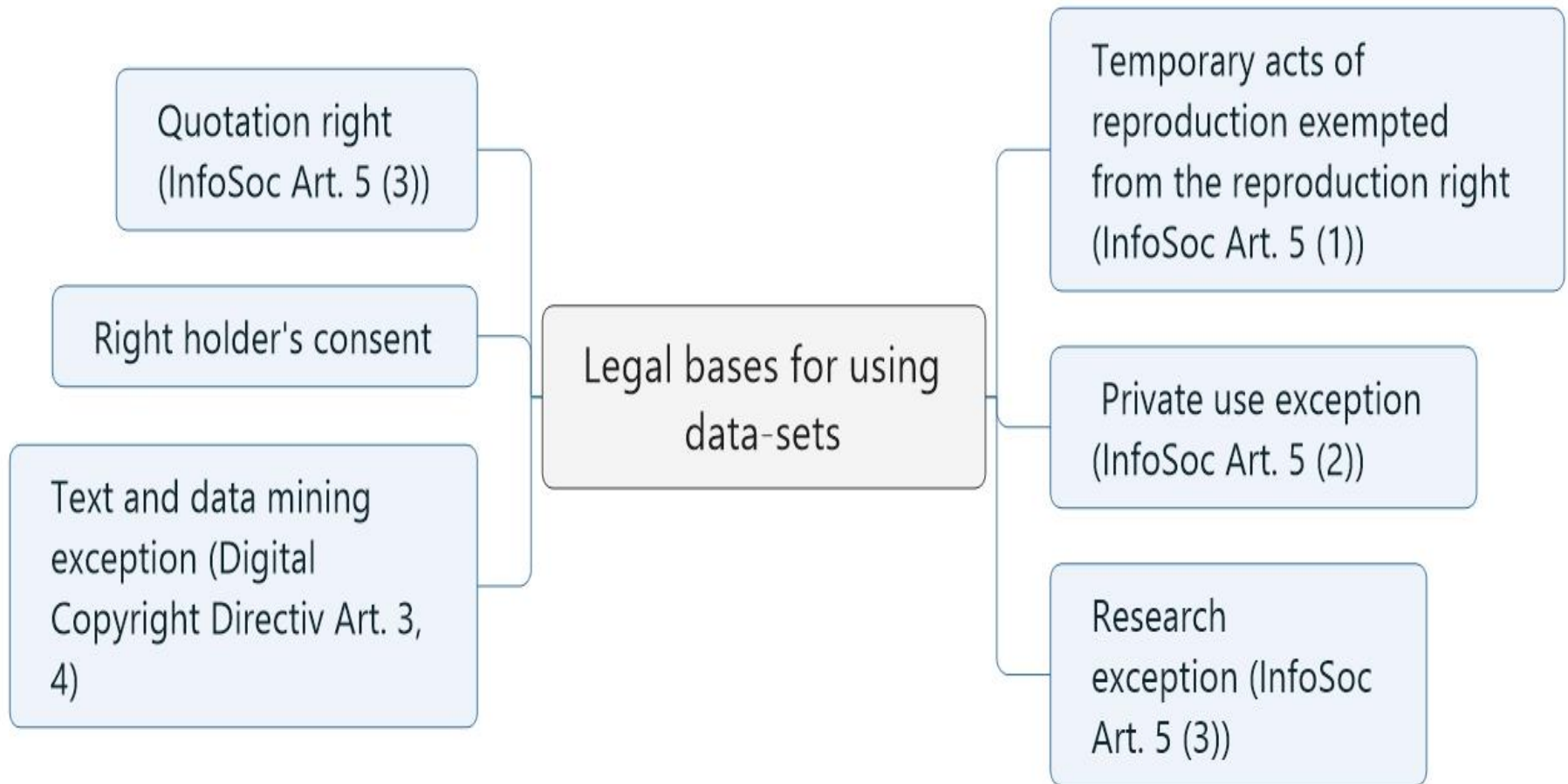
Language data & copyright:



Language data & personal data:



Use of LD & copyright protection:

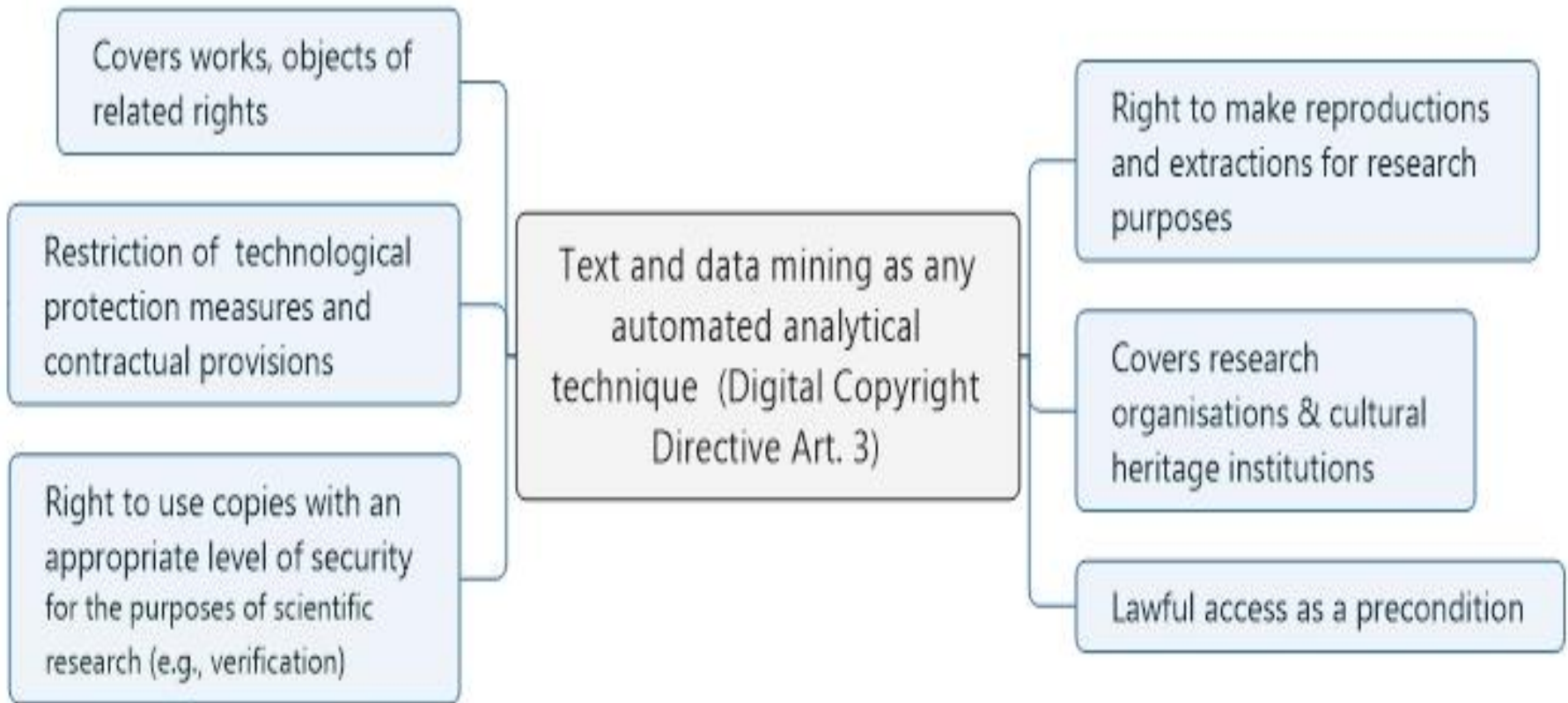


Use of LD & personal data protection:

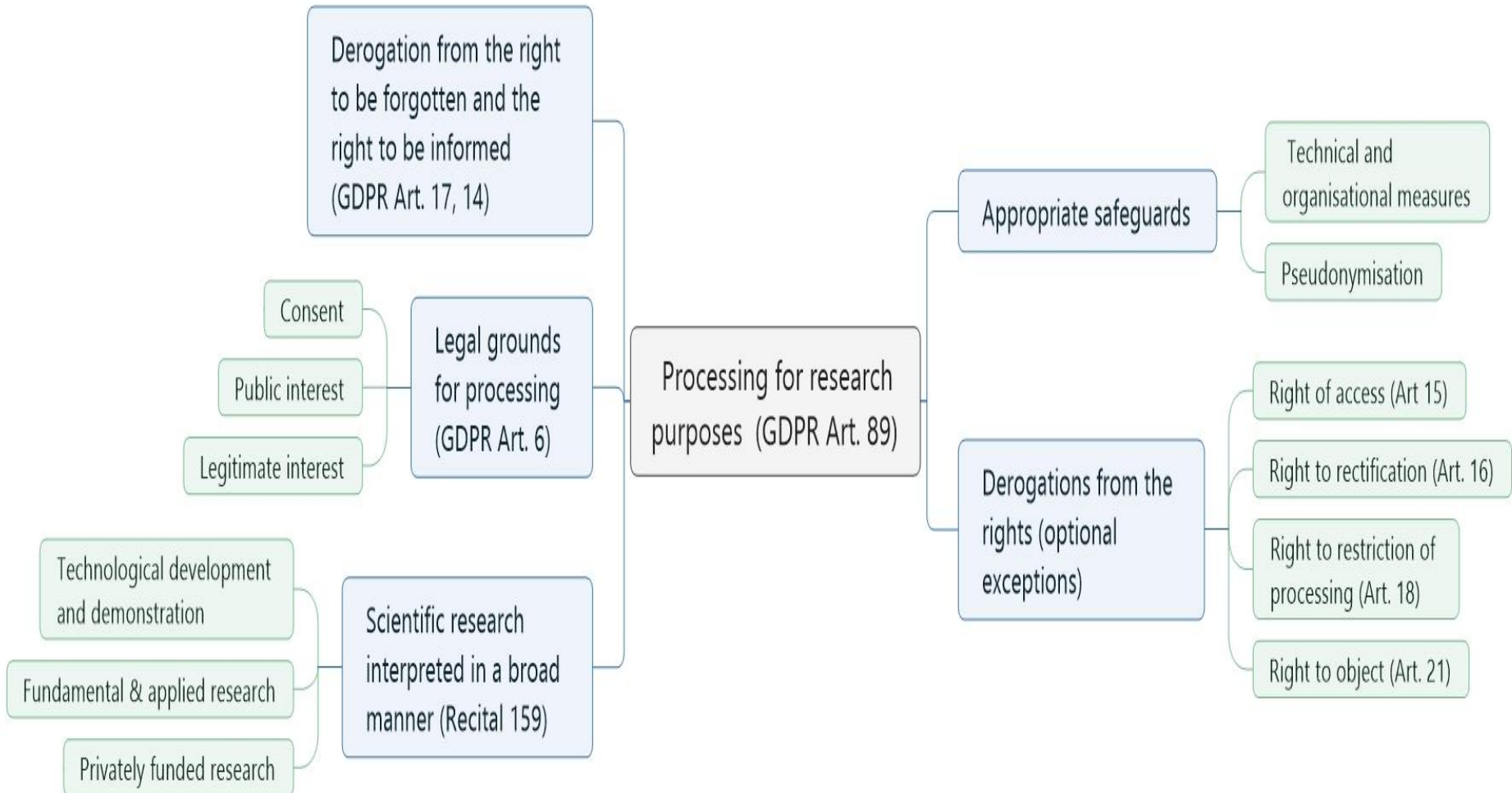




Research use & copyright:



Research use & personal data:

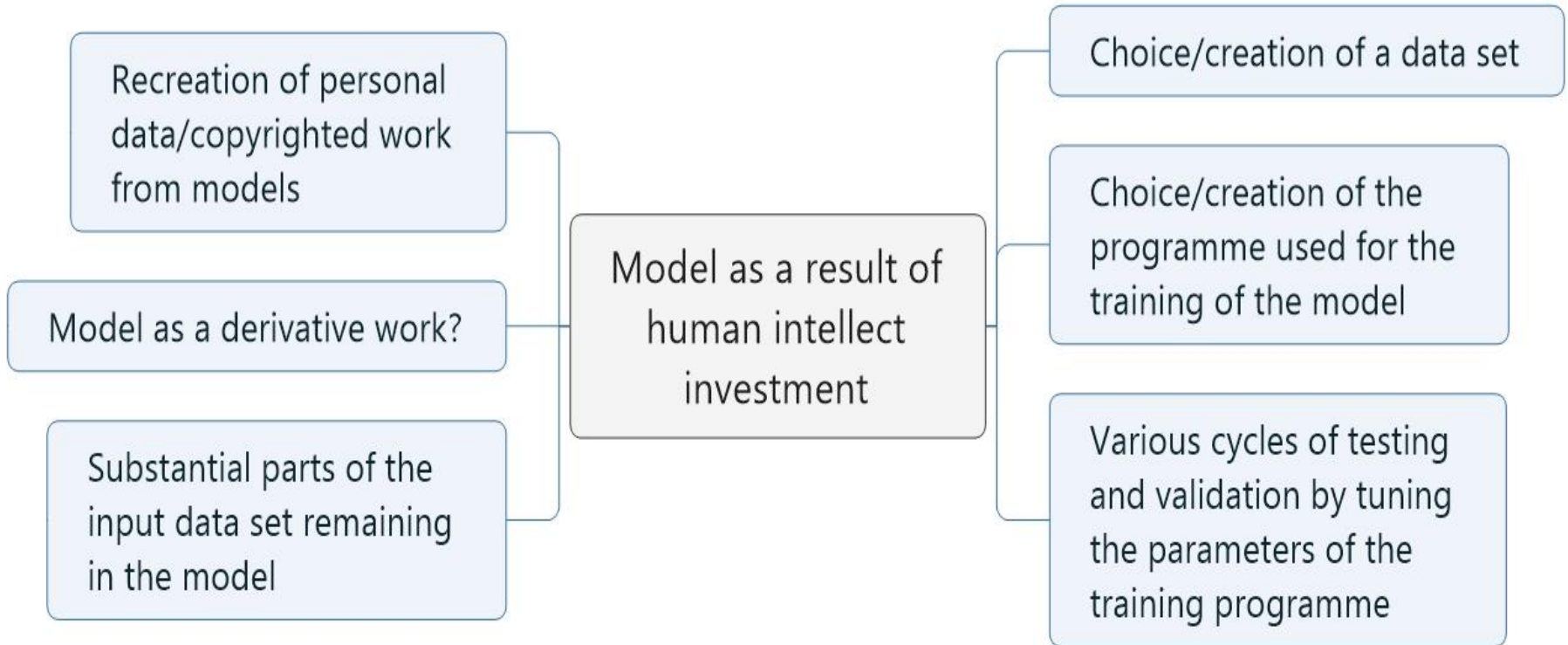


Models:



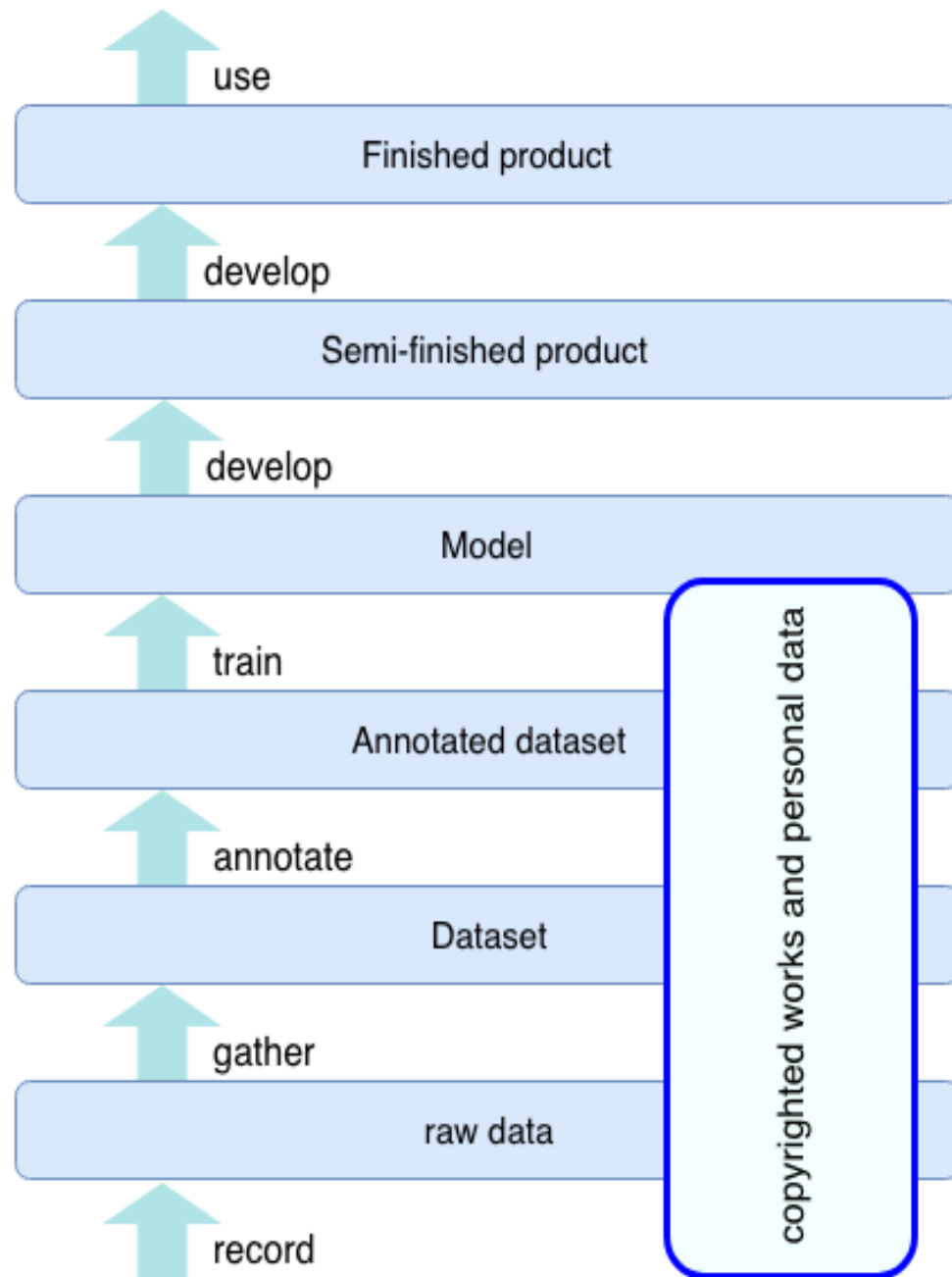


Models:

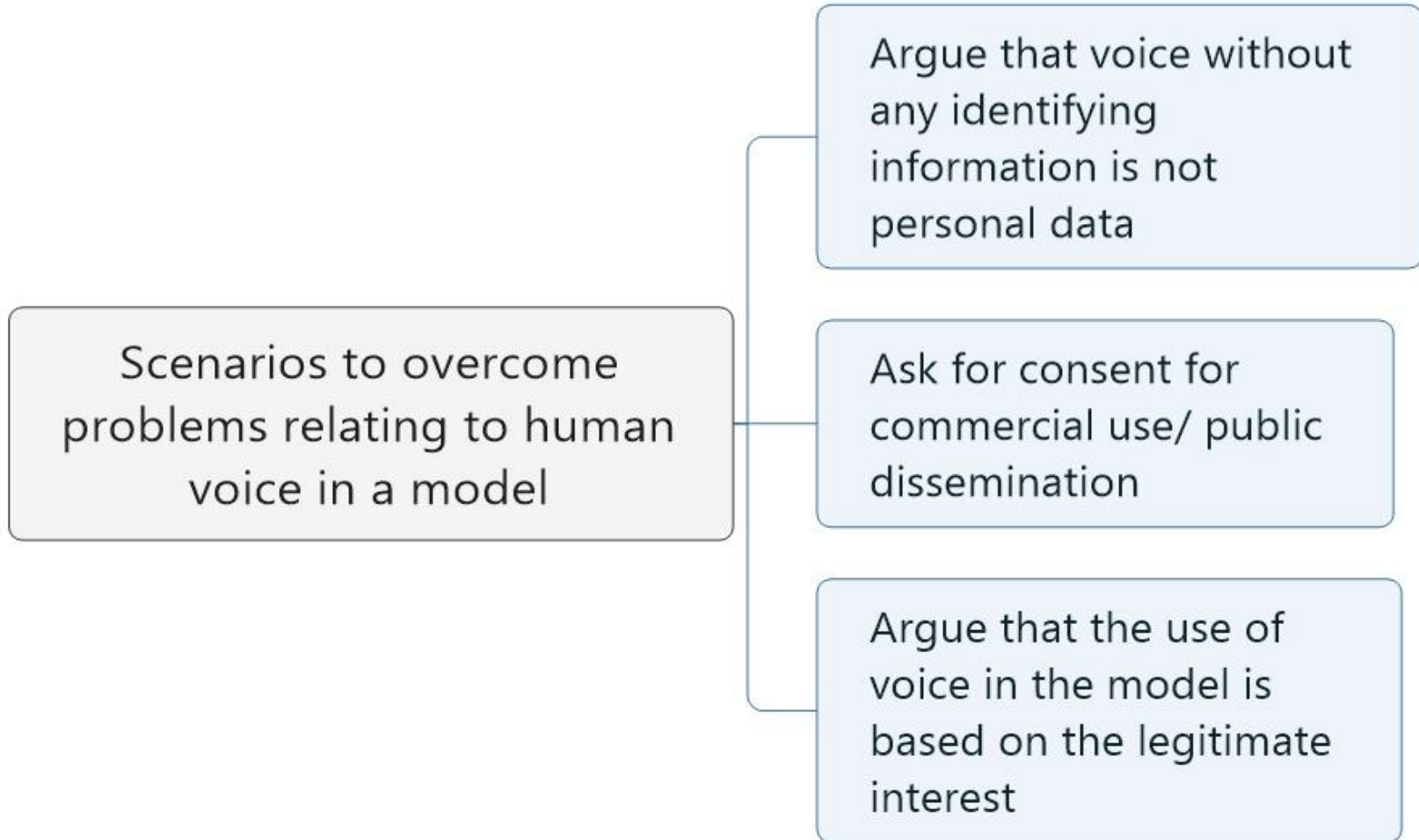




Models:



Personal data problems concerning voice:



Conclusions:

- Models are not necessarily subject to the same copyright and personal data restrictions as language data used as input;
- The limited impact of copyright and personal data restrictions on models potentially allows to:
 - Rely on favourable copyright/personal data provisions supporting research;
 - Make models publicly available;
 - Commercial use of models;
 - Personal data problems with voice need to be addressed.



TARTU ÜLIKOOL



Thank you