# Tackling Semantic :
## Semantic Analysis of Text Corpora using AI
### HU Berlin

**BIFOLD**

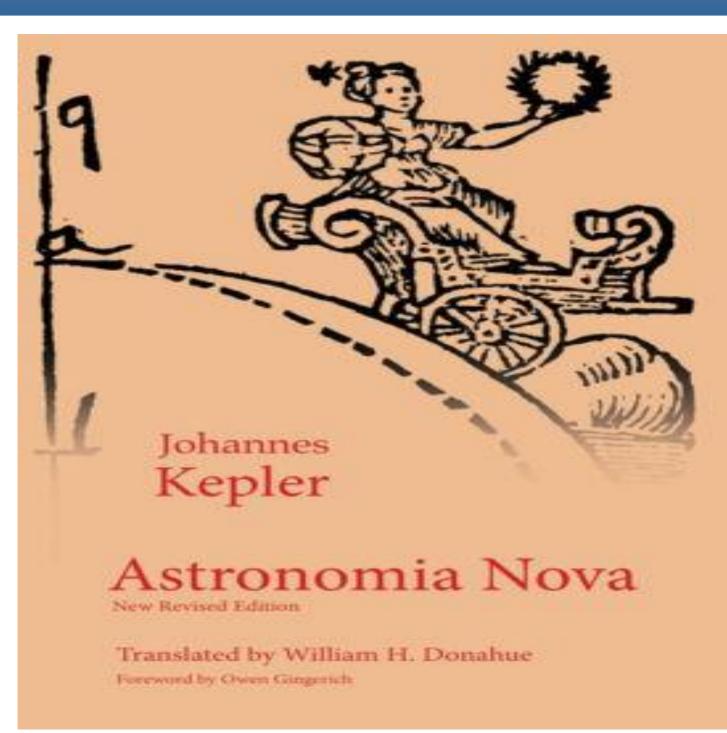## Analysis of Text Corpora

1) Kepler's book "Astronomia nova" Corpus
-
-
-

2) Exoplanet Publication Corpus [1]
-

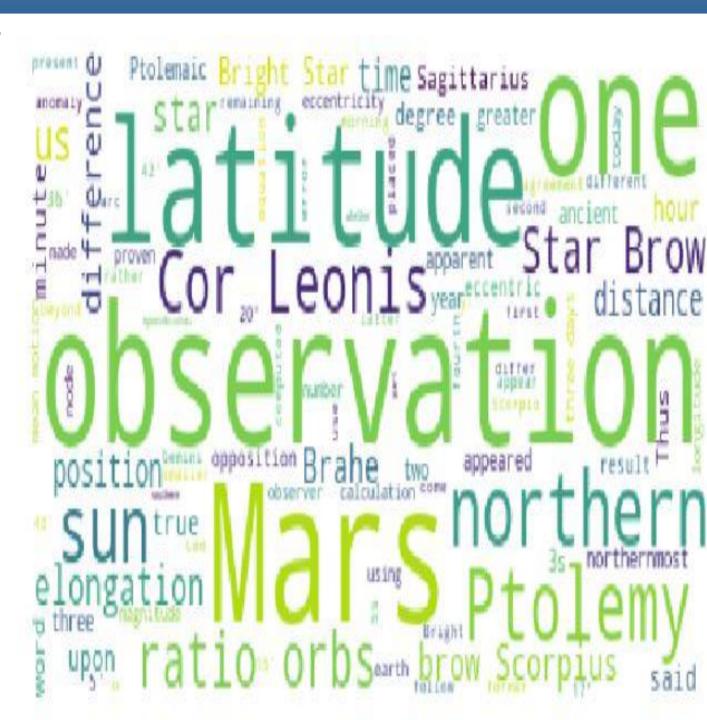3) Corona Virus Publication Corpus
-

## Structural Corpora

1) *Multi Index Dataframe*

2) *Sentence Structure with Spacy*
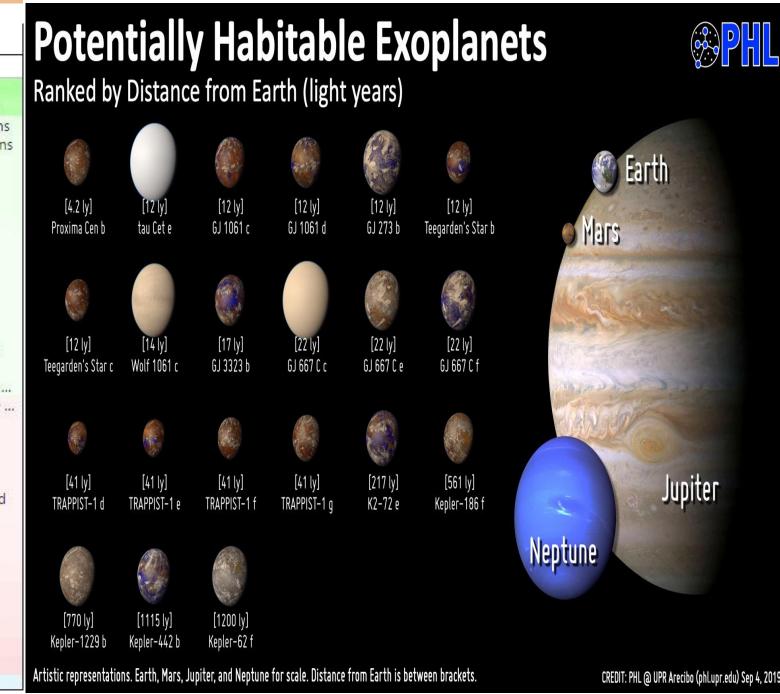
  a) Customized Entities

  b) Dependency Parsing

**Johannes Kepler — Astronomia Nova — New Revised Edition — Translated by William H. Donahue — Foreword by Owen Gingerich**

**Potentially Habitable Exoplanets** — Ranked by Distance from Earth (light years)

| y=0 top features | | y=1 top features | | y=2 top features | | y=3 top features | |
|---|---|---|---|---|---|---|---|
| Weight? | Feature | Weight? | Feature | Weight? | Feature | Weight? | Feature |
| +3.200 | present | +13.912 | compared | +15.586 | propose | +15.562 | compare |
| +2.823 | <BIAS> | +10.334 | compare | +9.788 | proposed | +12.624 | conclude |
| +1.953 | model | +2.055 | frequencies | +1.388 | macho | +2.129 | predictions |
| +1.565 | explore | +1.354 | derived | +1.354 | microlensing | +1.363 | conclusions |
| +1.321 | using | +1.589 | resulting | +1.196 | better | +1.305 | obtained |
| +1.161 | developed | +1.244 | oscillation | +1.126 | theta_ | +1.297 | eclipse |
| +1.102 | use | +1.147 | j_ | +0.999 | afta | +1.224 | existing |
| +0.920 | masses | +1.055 | specifically | +0.979 | blg | +1.213 | values |
| ... 6556 more positive ... | | +1.053 | simulations | +0.957 | series | +1.156 | depths |
| ... 1800 more negative ... | | +1.047 | trappist | +0.906 | enable | +0.978 | pressure |
| -0.925 | compares | +1.034 | previous | +0.900 | 35 | +0.906 | gyr |
| -0.938 | eclipse | +1.018 | axes | +0.872 | coronagraphic | +0.887 | opacity |
| -0.949 | significant | +1.015 | slow | +0.835 | eps | +0.880 | real |
| -0.955 | classical | +1.015 | curves | +0.818 | mission | +0.879 | event |
| -0.966 | observations | +1.014 | compares | +0.809 | interception | +0.872 | spectrum |
| -0.967 | obtained | +1.013 | methods | +0.783 | interactions | +0.871 | tio |
| -0.970 | clear | +0.998 | simpler | +0.761 | inner | +0.783 | inner |
| -0.997 | j_ | +0.997 | titan | +0.736 | libration | ... 7318 more positive ... | |
| -1.005 | theta_ | +0.994 | confirmed | ... 833 more positive ... | | -0.869 | planet |
| -1.054 | observed | +0.964 | case | ... 7523 more negative ... | | -0.887 | time |
| -1.055 | conclusions | ... 898 more positive ... | | -0.742 | develop | -0.910 | explore |
| -1.104 | macho | ... 7458 more negative ... | | -0.758 | transit | -0.949 | calculate |
| -1.159 | macho | -0.931 | eclipse | -0.832 | compare | -0.961 | developed |
| -1.166 | frequencies | -0.941 | investigate | -0.842 | temperature | -1.017 | previous |
| -1.197 | microlensing | -0.957 | modeling | -0.855 | surface | -1.104 | develop |
| -1.595 | better | -0.970 | transmission | -0.861 | effects | -1.112 | curves |
| -1.728 | predictions | -1.039 | atmospheric | -0.946 | evolution | -1.318 | apply |
| -8.344 | proposed | -1.045 | planet | -0.987 | mass | -1.935 | using |
| -11.004 | conclude | -1.097 | models | -1.063 | models | -1.970 | use |
| -11.936 | compared | -1.288 | present | -2.007 | present | -2.461 | present |
| -13.480 | propose | -1.436 | model | -2.208 | model | -2.855 | model |
| -19.334 | compare | -3.609 | <BIAS> | -3.589 | <BIAS> | -3.463 | <BIAS> |

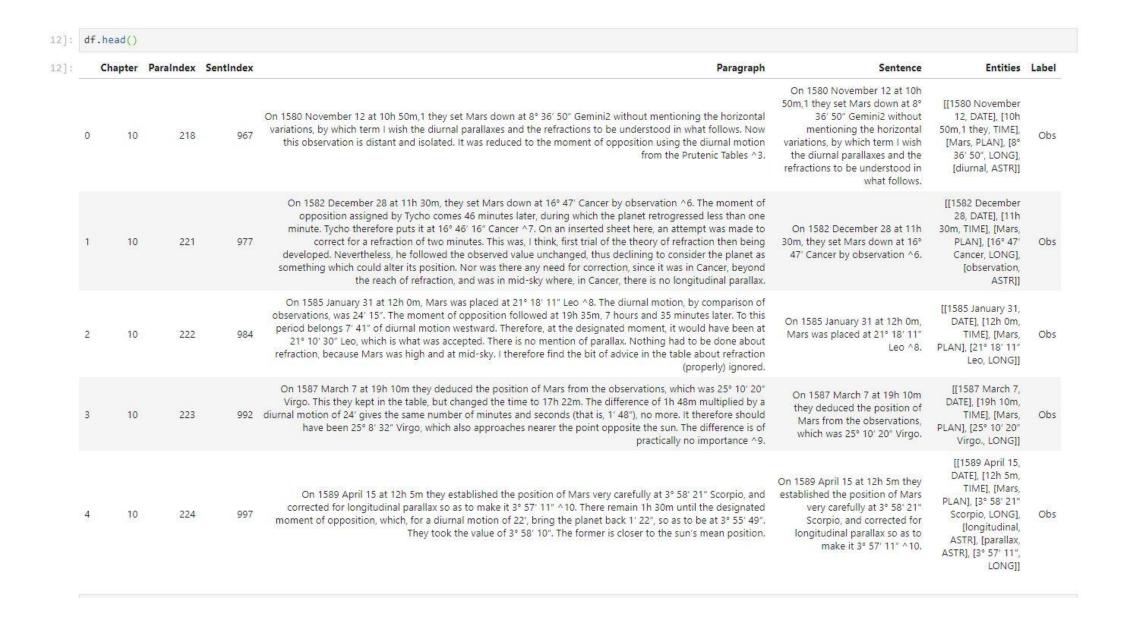## Customized Named Entity Recognition (NER) [2]

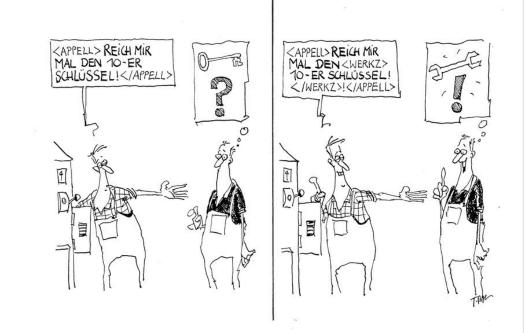*Key Idea.* Hybrid approach; Pattern based & deep learning to provide NER Model by developing spaCy.

## Explicit Observation Extraction [2]

*Key Idea.* Extract Semantic from the corpus using customized NER model.

['NAME', 'STAR', 'LONG', 'TIME', 'DATE', 'ASTR', 'GEOM', 'PARA', 'PLAN', 'LAT'] [421, 105, 1780, 208, 211, 4773, 606, 754, 934, 27]

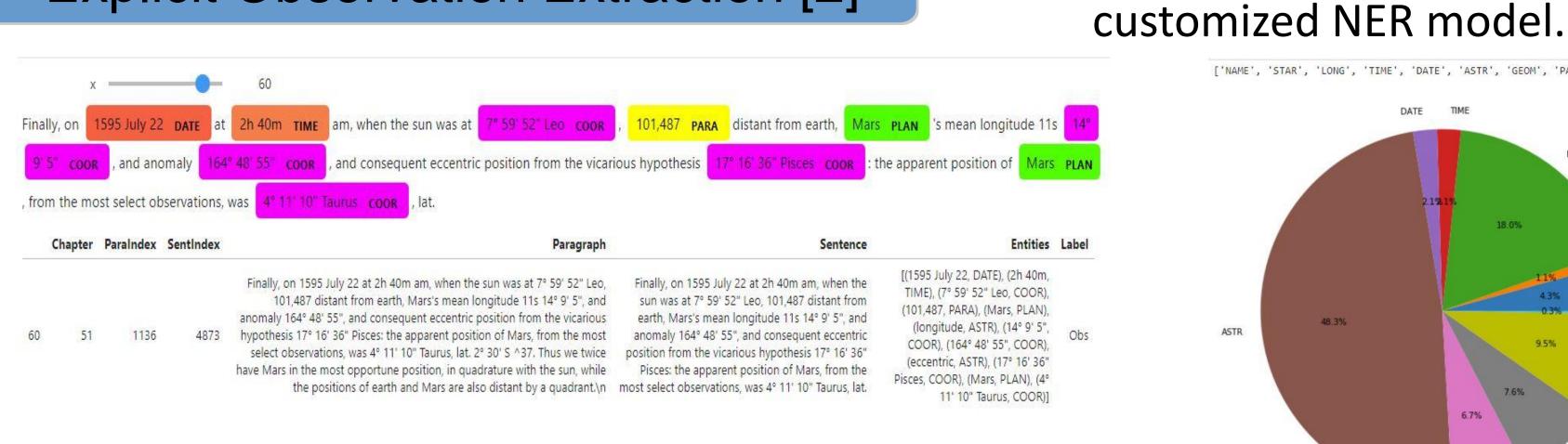## Future Direction: Semantic Objects

- *Tree structure*
  Idea. Integrate sentences using xml format to have semantic annotation for extracted observations as a first use case

- *Semantic Relation Extraction*
  Idea. Extract structure of argumentation components using semantic objects

## References

[1] Graßhoff, Gerd; Bier, Sabrina (2019): Database of abstracts in publications on exoplanets from the NASA archive. Zenodo. Dataset.
https://doi.org/10.5281/zenodo.3269732

[2] Yeghaneh, Mohammad; https://github.com/myeghaneh

[3] Ribeiro, M. T., Singh, S., & Guestrin, C. (2016, August). " Why should i trust you?" Explaining the predictions of any classifier. In *Proceedings of the 22nd ACM SIGKDD international conference on knowledge discovery and data mining* (pp. 1135-1144).

| | Sentence |
|---|---|
| 0 | Similarly, on 1591 January 22 at 7h in the morning, Mars was 34° 32' 45" from Spica with declination 17° 25' south, at an altitude of 10° 16. |
| 1 | Therefore, after correction for horizontal variations, the declination was 17° 30'. |
| 2 | Hence, the right ascension was 210° 22' 12", longitude 22° 33' Scorpio, latitude 1° 0' 30" north. |
| 3 | Now this time differs from ours by 1 day 19 hours, and the diurnal motion, from Magini, is 33'. |
| 4 | Therefore, 59' are required for the intervening time. |