



Vers une évaluation universelle du niveau de complexité des textes

Marc Benzahra

Superviseur **François Yvon**

Equipe **TLP**

Bureau **A203**

Financement **Glose**

Text complexity definition

Readability is what makes some **texts easier to read than others** at the **content** aspect of a text

e.g. Word length, sentence length, number of syllables, paragraph length, words POS tags and dependency tree

Whereas **lisibility** is about the **form** of a text

e.g. Words spacing, text font, character size

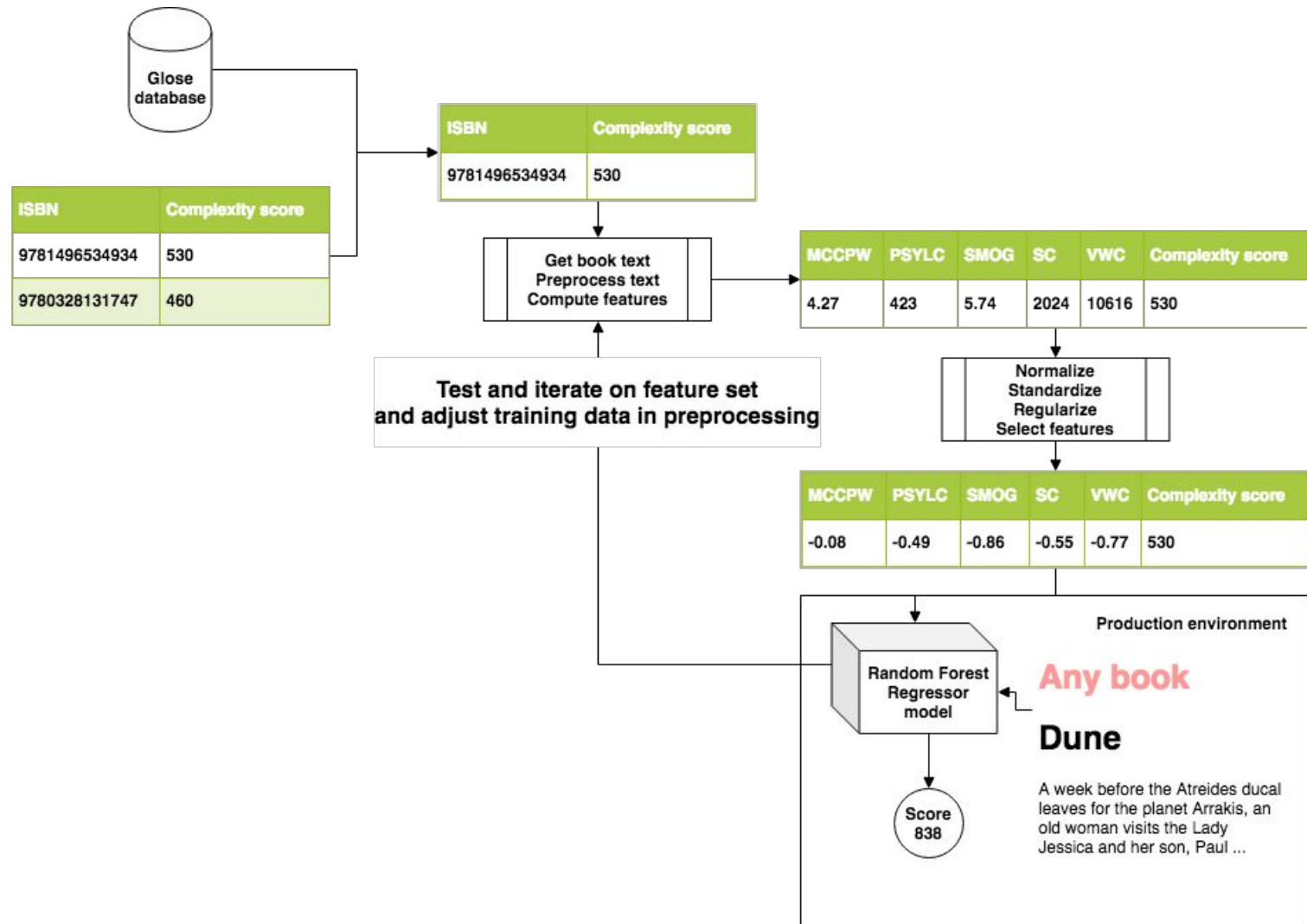
Text complexity usage

Convey information to most readers

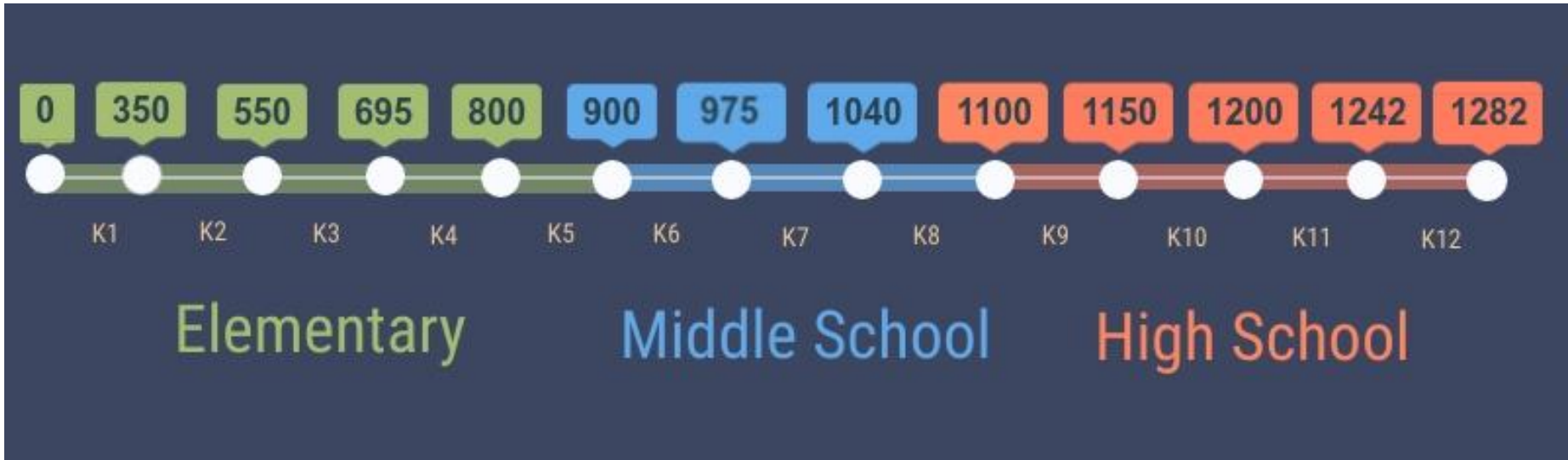
e.g. Drugs leaflets, news information, administrative and legal texts

Engage readers through content that is **gradually more difficult**

e.g. First and second language learners, ideally to any domain learning condition



Classroom levels scale



Novel approach: unsupervised task

Hypothesis

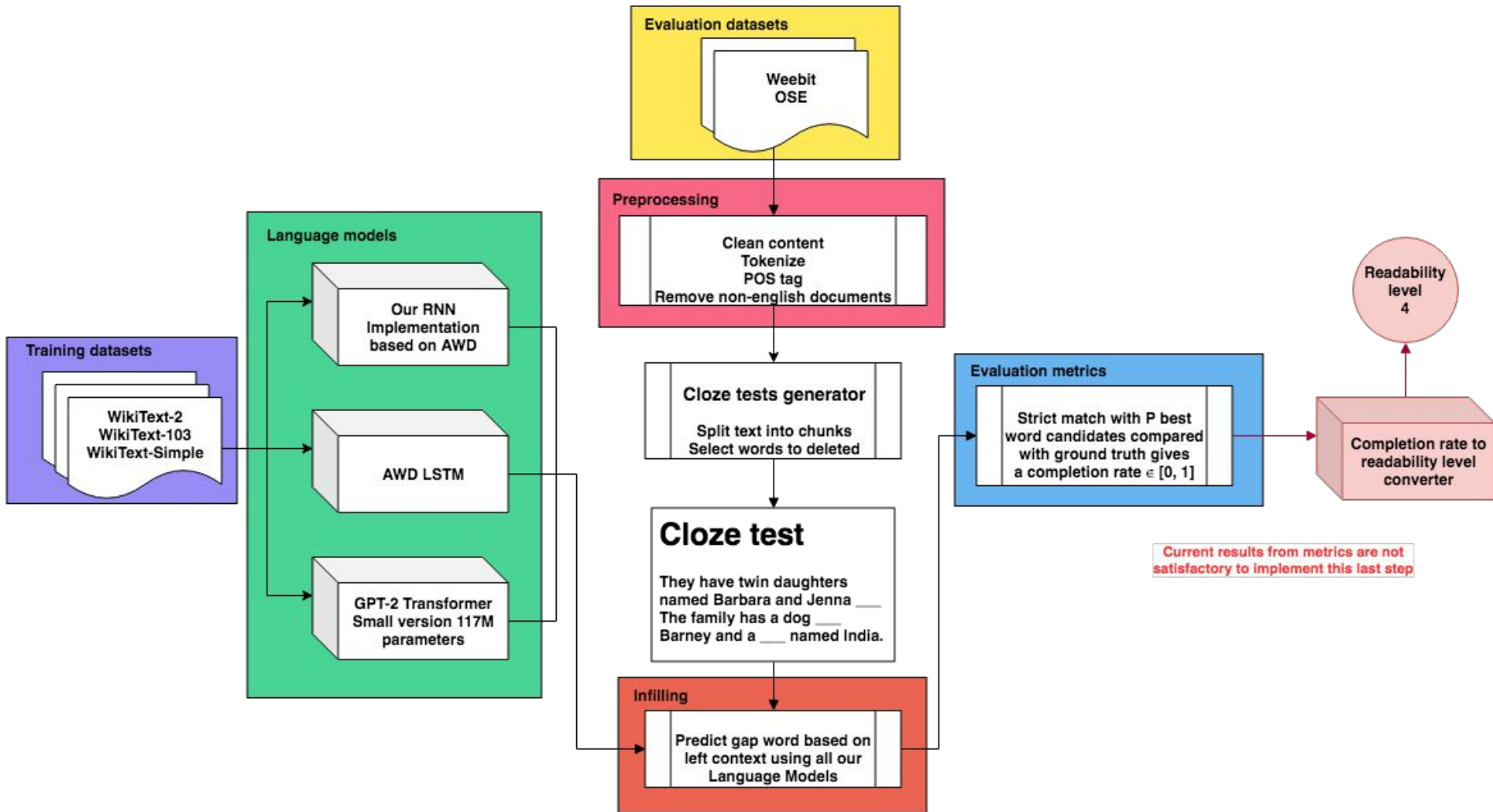
The simpler a text is, the better it should be understood by a machine

Idea

We should expect a strong correlation between readability levels and performance of language models at infilling Cloze tests

Current state

Correlations observed with our systems are still really small compared to other readability metrics such as Flesch-Kincaid or the previous model



Merci pour votre attention



**« Créer, c'est ainsi donner
une forme à son destin. »**

Albert Camus, *Le Mythe de Sisyphe*

References

William H Dubay, 2007

Smart Language

In Readers, Readability, and the Grading of Text. Costa Mesa: Impact Information.

Thomas François and **Cédrick Fairon**, 2012

An AI readability formula for French as a foreign language

In Proceedings of the 2012 Joint Conference on Empirical Methods in NATural Language Processing and Computational Natural Language Learning, pages 466-477.

Stephen Merity, **Nitish Shrirish Keskar** and **Richard Socher**, 2018

Regularizing and optimizing LSTM language models

In Proceedings of the International Conference on Learning Representations, ICLR.

References

Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei and Ilya Sutskever, 2019

Language models are unsupervised multitask learners

Technical report, OpenAI.

Sarah E. Petersen and Mari Ostendorf, 2009

A machine learning approach to reading level assessment

Computer Speech & Language 23(1):89 - 106.

Sowmya Vajjala and Ivana Lucic, 2018

Onestopenglish corpus: A new corpus for automatic readability assessment and text simplification

In Proceedings of the Thirteenth Workshop on Innovative Use of NLP for Building Education Applications, pages 297-304, New Orleans, Louisiana.