N-gram analysis of amyloid data

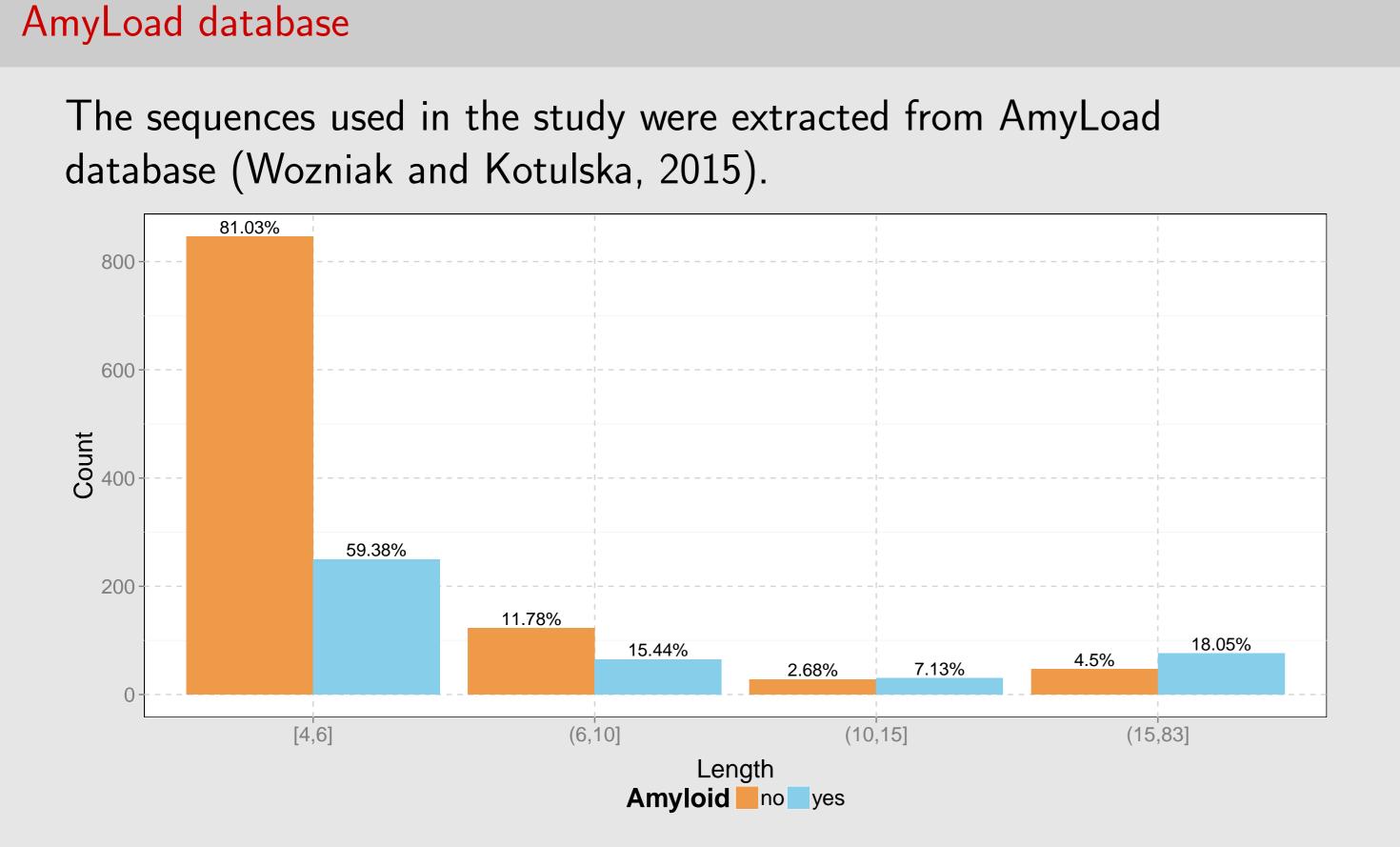
 $\sf Michał Burdukiewicz^1$, $\sf Piotr Sobczyk^2$, $\sf Paweł Mackiewicz^1$ and $\sf Małgorzata Kotulska^3$

¹University of Wrocław, Department of Genomics, Poland

²Wrocław University of Technology, Department of Mathematics, Poland

³Wrocław University of Technology, Department of Biomedical Engineering, Poland

Introduction	More stuff
Stuff	Stuff



The height of the bar is equal to the number of sequences within given length interval in the database. The percentage above the bar represent fraction sequences within given length interval among amyloid and non-amyloid sequences.

Bibliography

Kawashima, S., Pokarowski, P., Pokarowska, M., Kolinski, A., Katayama, T., and Kanehisa, M. (2008). AAindex: amino acid index database, progress report 2008. *Nucleic Acids Research*, 36(suppl 1):D202–D205.

Wozniak, P. P. and Kotulska, M. (2014). Characteristics of protein residue-residue contacts and their application in contact prediction. *Journal of Molecular*

Modeling, P. P. and Kotulska, M. (2014). Characteristics of protein residue-residue contacts and their application in contact prediction. Journal of Nombre Modeling, 20(11).

Wozniak, P. P. and Kotulska, M. (2015). AmyLoad: website dedicated to amyloidogenic protein fragments. Bioinformatics (Oxford, England).

Clustering of amino acids

contactivity (Wozniak and Kotulska, 2014)

- 1. 10 physicochemical scale were extracted from AAIndex database (Kawashima et al., 2008) and normalized between normalized between 0 i 1.
- 2. All combinations of characteristics (only one scale per property) were clustered using Euclidean distance and Ward's method.
- 3. Each clustering was divided into 3 to 6 groups.