# Methods for Text Complexity Assessment in Russian

### RuSTEP September 2022

Vladimir Ivanov, Innopolis University



### **Text complexity Definition / Multilevel Phenomenon**

- TC has been studied at various levels of linguistic units
- whole texts
- sentence level (Schumacher et al. (2016); lavarone et al. (2021)
- individual words (Shardlow et al. (2021, 2020))

(Crossley et al. (2008); Collins-Thompson and Callan (2005); Heilman et al. (2008))

## **Text Complexity as Readability Ease Readability indices / formulas**

readability formulas tools to match texts and readers

Score	School level (US)	
100.00–90.00	5th grade	Very easy to r
90.0–80.0	6th grade	Easy to read.
80.0–70.0	7th grade	Fairly easy to
70.0–60.0	8th & 9th grade	Plain English.
60.0–50.0	10th to 12th grade	Fairly difficult
50.0–30.0	College	Difficult to rea
30.0–10.0	College graduate	Very difficult to
10.0–0.0	Professional	Extremely diff

### Notes

read. Easily understood by an average 11-year-old student.

Conversational English for consumers.

read.

Easily understood by 13- to 15-year-old students.

to read.

۱d.

o read. Best understood by university graduates.

ficult to read. Best understood by university graduates.

# FK reading ease

Average syllables per word



https://en.wikipedia.org/wiki/Flesch-Kincaid\_readability\_tests



# Applications

0





## **Flesch-Kincaid** 1975

- Flesch-Kincaid (English)

- Adopted by I.Oboroneva (2006), for Russian

# The Flesch Reading Ease (FRE), and the Flesch–Kincaid Grade Level (FKG) $FRE = 206.835 - 1.015 \times ASL - 84.6 \times AWS$ $FKG = 0.39 \times ASL + 11.8 \times AWS - 15.59$

 $FRE = 206.835 - 1.52 \times ASL - 65.14 \times AWS$ 

# **Other indices / readability formulae**

- Gunning Fog Index
- Coleman Liau Index
- SMOG Index

. . .

• Dale-Chall Readability Formula

https://readable.com/features/readability-formulas/

### Modified version of FKG for Russian Solovyev, Solnyshkina, Ivanov, Batyrshin (2018)

• Linear regession fitted on School Textbooks

 $FKG = 0.36 \times ASL + 5.76 \times AWS - 11.97$ 

# Studying of Groups of Features

An extended feature set for the text explored:

- Features based on length and frequency of words and sentences
- Features based on Part-of-Speech tags
- Features based of syntactic dependencies

# Features based on length and frequency

- •ASL is an average number of words per sentence
- •ASW is an average number of syllables per word
- •FREQ is a cumulative frequency of content words

# Features based on POS tags

- •NOUNS is a number of nouns per sentence
- •**VERBS** is a number of verbs per sentence
- •**ADJ** is a number of adjectives per sentence
- •**PRONOUNS** is a number of pronouns per sentence
- •PERONAL PRONOUNS is a number of personal pronouns per sentence •NEG is a number of negations per sentence

# Features based of syntactic dependencies

**1.AVERAGE\_PATH** is the quotient of the number of nodes and the number of leaves in a sentence 2.AVERAGE\_SOCHIN\_LENGTH is the average length of coordinating constructions **3.DEEPRICH\_RATE** is the average number of verbal participles **4.DEEPRICH\_V** is the average span of a verbal adverb phrase **5.LEAVES\_NUMBER** is the average number of 'leaves' in a sentence **6.LONGEST\_PATH** is the average length of the longest branch links

**9.PODCHIN\_RATE** is the average number of subordinate links that has at least one dependent

the participle

**12.SENTSOCH\_NUMBER** is the average number of compound sentences **13.SOCHIN\_NUMBER** is defined as the average number of coordinating chains **14.PATH\_NUMBER** is defined as the average number of sub-trees (in a sentence) ignored.

- **7.NOUNS\_DEP** is the average number of modifiers in a nominal group; coordinating and explanatory links are ignored **8.PODCHIN\_NUMBER** is the ratio of sentences in which there is at least one subordinate conjunctions or relational
- **10.PRICH\_RATE** is the average number of participial construction; participial constructions are defined as a participle
- **11.PRICH\_V** is the average span of a participial construction is the quotient of the number of nodes that depend on
- **15.VERBS\_DEP** is defined as the average number of finite dependent verbs and is calculated as the sum of nodes directly dependent on the finite verb divided by the number of finite verbs; coordinating and explanatory links were

# Correlation between features and grade level

	Feature name	Correlation		Feature name	Correlation
1	ASL	0.94	13	NOUNS	0.82
2	ASW	0.94	14	VERBS	0.74
3	SOCHIN_NUMBER	0.93	15	NEGATIONS	0.7
4	PRICH_RATE	0.91	16	PRONOUNS	0.7
5	NOUNS_DEP	0.88	17	PODCHIN_RATE	0.64
6	AVERAGE_SOCHIN_LEN	0.87	18	PODCHIN_NUMBER	0.62
7	PATH_NUMBER	0.87	19	DEEPRICH_V	0.52
8	LONGEST_PATH	0.84	20	PERS_PRONOUNS	0.47
9	FREQ	0.84	21	DEEPRICH_RATE	0.44
10	LEAVES_NUMBER	0.84	22	VERBS_DEP	0.43
11	AVERAGE_PATH	0.84	23	PRICH_V	0.33
12	ADJ	0.82	24	SENTSOCH_NUMBER	0.03

# Significance of features

	Feature	Absolute value of Coefficient in Ridge Regression
1	ASL	0.506
2	ASW	0.125
3	SOCHIN_NUMBER	0.119
4	PRICH_RATE	0.106
5	LONGEST_PATH	0.089
6	PATH_NUMBER	0.079
7	LEAVES_NUMBER	0.075
8	AVERAGE_SOCHIN_LE	0.071
9	NOUNS_DEP	0.071
10	FREQ	0.034
11	NEGATIONS	0.01
12	AVERAGE_PATH	0.007
13	PERS_PRONOUNS	0.003
14	VERBS	0.001
15	ADJ	0.001
16	NOUNS	14 0.0

## Study of fragment size V. Solovyev, V. Ivanov, M. Solnyshkina (2017)



# Study of fragment size



	1-3	399 86%	43 11%	106 6%	8 0%
put Class	4-6	57 12%	273 71%	202 11%	25 1%
Out	7-9	3 1%	56 14%	782 44%	339 8%
1	0-11	0 0%	12 3%	668 38%	3628 91%
		1-3	4-6	7-9	10-11

Target Class

**Beyond Readability formulae** 

## Analysis of Academic text Complexity in collaboration with V. D. Solovyev and M.I. Solnyshkina (KFU)

- Classic ML (linear regression, classifiers)
  - 2018: Analysis of sets of features and fragment size
  - 2019: Analysis of semantic-level features
- Deep learning
  - 2020-21: Appication of neural networks to texts
  - 2021-22: Complexity of sentences and words

# **Corpora of Academic texts**

- NIK + BOG Corpus (14 books)
- + Textbooks in History (approx. 5 books)
- + Textbooks in Biology (approx. 5 books)
- + Elementary school (>100 books)

tokenization, splitting text into sentences, excluded all extremely long sentences and short sentences

Cue de laval	Tok	ens	Sentences		ASL		
Grade level	BOG	NIK	BOG	NIK	BOG	NIK	B
5-th	-	17 221	-	1 499	-	11.49	
6-th	16 467	16 475	1 273	1 197	12.94	13.76	2.
7-th	23 069	22 924	1 671	1 675	13.81	13.69	2.
8-th	49 796	40 053	3 181	2 889	15.65	13.86	2.
9-th	42 305	43 404	2 584	2 792	16.37	15.55	3.
10-th	75 182	39 183	4 468	2 468	16.83	15.88	3.
10-th*	98 034	-	5 798	-	16.91	-	3.
11-th	-	38 869	-	2 270	-	17.12	
11-th*	100 800	-	6 004	-	16.79	-	3.



### Corpora of Academic texts File M Sentence level

- Dataset of sentence pairs
- Dataset of sentences



File Name (book)	Grade Level	Label	Category	ASL	AWL	Total S
year_bog_11p.txt	11.50	11	high	18.63	7.13	
year_petrov_11.txt	11.00	11	high	15.65	6.71	
year_guryan_11.txt	11.00	11	high	16.01	6.72	
year_nik_11.txt	11.00	11	high	17.99	6.95	
year_ponom_11.txt	11.00	11	high	15.71	6.65	
year_plenko_11.txt	11.00	11	high	16.66	6.63	
year_bog_10p.txt	10.50	11	high	19.07	6.92	
year_sobol_10.txt	10.00	10	high	15.93	6.75	
year_unk_10.txt	10.00	10	high	16.19	6.68	
year_nik_10.txt	10.00	10	high	17.81	6.91	
year_bog_10.txt	10.00	10	high	18.27	6.78	
year_klimov_10.txt	10.00	10	high	17.09	6.76	
year_bog_09.txt	9.00	9	high	17.88	6.68	
year_nik_09.txt	9.00	9	high	16.90	6.79	
year_bog_08.txt	8.00	8	medium	17.49	6.72	
year_nik_08.txt	8.00	8	medium	15.74	6.41	
year_nik_07.txt	7.00	7	medium	15.41	6.14	
year_bog_07.txt	7.00	7	medium	15.00	6.46	
year_nik_06.txt	6.00	6	medium	15.94	6.18	
year_bog_06.txt	6.00	6	medium	15.13	5.86	
year_nik_05.txt	5.00	5	medium	13.11	5.57	
year_vah_4pu.txt	4.50	5	medium	15.78	5.86	
year_vah_4u.txt	4.00	4	low	13.78	6.17	
year_ben_4u.txt	4.00	4	low	12.72	6.38	
year_gor_4u.txt	4.00	4	low	14.75	6.56	
year_rud_3u.txt	3.00	3	low	14.15	5.67	
year_vah_2pu.txt	2.50	3	low	11.13	5.75	
year_uch_2pu.txt	2.00	2	low	14.14	5.73	
year_uch_2u.txt	2.00	2	low	12.00	6.05	
year_vah_2u.txt	2.00	2	low	11.10	5.73	
yead_rud_2u.txt	2.00	2	low	13.40	5.44	
year_vah_1pu.txt	1.50	2	low	11.22	5.69	
year_rag_1u.txt	1.00	1	low	8.76	5.95	
year_rog_1u.txt	1.00	1	low	10.33	5.95	
year_rud_1u.txt	1.00	1	low	12.20	5.17	
year_lut_1u.txt	1.00	1	low	9.74	6.36	
year_kur_1u.txt	1.00	1	low	9.86	6.15	
vear vah 10.txt	1.00	1	low	11.17	5.38	



5,648 5,029 5,848 2,078 2,190 3,470 5,579 5,236 2,000 2,271 3,145 3,967 1,710 2,480 2,999 1,821 1,509 1,632 1,029 985 1,566 1,174 1,423 604 833 1,319 1,005 1,559 1,621 1,100 619 292 74 468 495 390 200 139

## Tasks and model types depending on the dataset we have different task setups

Dataset	Regression	Binary Classification	Multiclass Classification
1-s dataset	grade level (value)	-	11 or 3 categories
2-s dataset	difference between grade levels	complex / simple sentence	

### **Types of Models:**

- Linear regression
- Transformer-based
- GNN-based

### **Model architectures** for single sentence complexity prediction





### Tokens of a sentence

### **Model architectures** for classification of pairs of sentences



SVM Classification (binary)

### **Model architectures Graph Neural Networks**



### Multi-layer Graph Convolutional Network (GCN) applied for Dependency Tree



fastText embeddings

### **Results** Regression

Model

Linear Reg. on 2 parameters (SL+AWL) Linear Reg. on Sentence Embeddings SVR on Sentence Embeddings Fine-tuned RuBERT GNN

	1-s data			2-s data	
$R^2$	MSE	MAE	$\mid R^2$	MSE	MAE
0.13 0.65 0.71 0.80	8.34 3.39 2.79 1.96	2.32 1.37 1.11 0.80	- 0.73 0.77 0.98	- 10.55 8.96 0.79	2.56 2.28 0.56

### **Results** Classification

2-s dataset (with margin=3)

Model

SVM on Sentence Embeddings Fine-tuned RuBERT GNN

Accuracy	F1	Ρ	R
0.93	93.18	93.18	93.18
0.98	98.47	98.44	98.50
0.97	96.60	96.60	96.60

## **Complexity of words** Lexical complexity prediction

- Commonly this task is referred to as Complex Word Identification (CWI) or Lexical Complexity Prediction (LCP).
- Following (Paetzold, 2016; Zampieri 2017; Yimam 2018; Shardlow, 2021)
- Results for Russian (Abramov, Ivanov, 2022)
  - a corpus consisting of 931 distinct words that occurred within 3,364 different contexts
  - We evaluate a linear regression model as a baseline
  - handcrafted features, fastText and ELMo embeddings of target words.

## **Annotation Results** from Holy Bible







# **Modeling Results**

Table 2: The result of linear regression on handcrafted features (HC), Fasttext and ELMo

	Handcrafted	Fasttext	ELMo	Fasttext+HC	ELMo+HC
MAE	0.102	0.084	0.099	0.084	0.099
Pearson correlation	0.342	0.614	0.498	0.619	0.501

	Frequency	Word length	Number of syllables
Frequency	1	-0.206	-0.172
Word length	-0.206	1	0.819
Number of syllables	-0.172	0.819	1

embeddings and concatenated features

Table 3: Pearson correlation between handcrafted features

# Summary

- Text complexity is a multilayer phenomenon
  - whole text, passage, sentence, word
  - features
- Domain / genre dependent  $\bullet$
- Multilingual text complexity analysis

frequency / statistical features, syntactical features, semantic, contextual

## References

- 409-425.
- Solovyev, Valery, et al. "Prediction of reading difficulty in Russian academic texts." Journal of intelligent & fuzzy systems 36.5 (2019): 4553-4563.
- Cham, 2018.
- Tools and Resources to Empower People with REAding Difficulties (READI). 57-62.
- Semantic Evaluation (SemEval-2021). 1-16. DOI: 10.18653/v1/2021.semeval-1.1
- Association for Computational Linguistics), 186–199. doi:10.18653/v1/2021.cmcl-1.23
- 1462
- Crossley, S. A., Greenfield, J., and McNamara, D. S. (2008). Assessing text readability using cognitively based indices. Tesol Quarterly 42, 475–493
- innovative use of NLP for building educational applications. 71–79
- 560-569. DOI: 10.18653/v1/S16-1085
- International Workshop on Semantic Evaluation (SemEval-2016). 1001-1005. DOI: 10.18653/v1/S16-1155
- Workshop on Natural Language Processing Techniques for Educational Applications (NLPTEA 2017). 59-63.

Abramov, Aleksei V., and Vladimir V. Ivanov. "Collection and evaluation of lexical complexity data for Russian language using crowdsourcing." Russian Journal of Linguistics 26.2 (2022):

Solnyshkina, Marina, Vladimir Ivanov, and Valery Solovyev. "Readability formula for Russian texts: a modified version." Mexican International Conference on Artificial Intelligence. Springer,

Shardlow, Matthew, Michael Cooper & Marcos Zampieri. 2020. CompLex — A New Corpus for Lexical Complexity Prediction from Likert Scale Data. Proceedings of the 1st Workshop on

Shardlow, Matthew, Richard Evans, Gustavo H. Paetzold & Marcos Zampieri. 2021. Semeval-2021 task 1: Lexical complexity prediction. Proceedings of the 15th International Workshop on

lavarone, B., Brunato, D., and Dell'Orletta, F. (2021). Sentence complexity in context. In Proceedings of the Workshop on Cognitive Modeling and Computational Linguistics (Online:

Collins-Thompson, K. and Callan, J. (2005). Predicting reading difficulty with statistical language models. Journal of the american society for information science and technology 56, 1448-

Schumacher, E., Eskenazi, M., Frishkoff, G., and Collins-Thompson, K. (2016). Predicting the relative difficulty of single sentences with and without surrounding context. In Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing (Austin, Texas: Association for Computational Linguistics), 1871–1881. doi:10.18653/v1/D16-1192

Heilman, M., Collins-Thompson, K., and Eskenazi, M. (2008). An analysis of statistical models and features for reading difficulty prediction. In Proceedings of the third workshop on

Paetzold, Gustavo & Lucia Specia. 2016. SemEval 2016 Task 11: Complex Word Identification. In Proceedings of the 10th International Workshop on Semantic Evaluation (SemEval-2016).

Zampieri, Marcos, Liling Tan & Josef van Genabith. 2016. Macsaar at semeval-2016 task 11: Zipfian and character features for complex word identification. In Proceedings of the 10th

Zampieri, Marcos, Shervin Malmasi, Gustavo Paetzold & Lucia Specia. 2017. Complex Word Identification: Challenges in Data Annotation and System Performance. Proceedings of the 4th