

LT3 & SDL

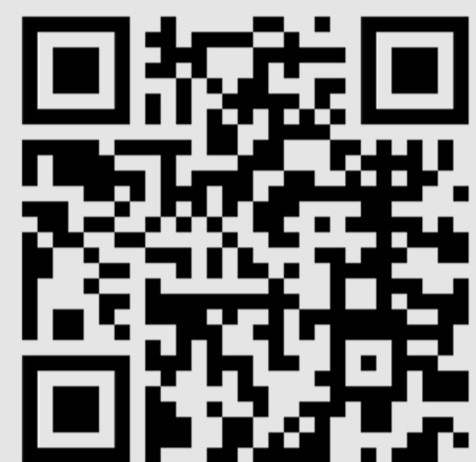
Orphée De Clercq, Joni Kruijsbergen, Fauve De Backer & Goedele Vandommele

# HOW TO RELIABLY ASSESS DUTCH WRITING SKILLS: MAN VS MACHINE

2024 = LAUNCH OF CENTRALIZED TESTING IN FLANDERS

STEUNPUNT  
CENTRALE TOETSEN  
IN ONDERWIJS

Ca. 70,000 pupils  
MATHEMATICS & DUTCH  
Multidisciplinary team



DUTCH SKILLS?



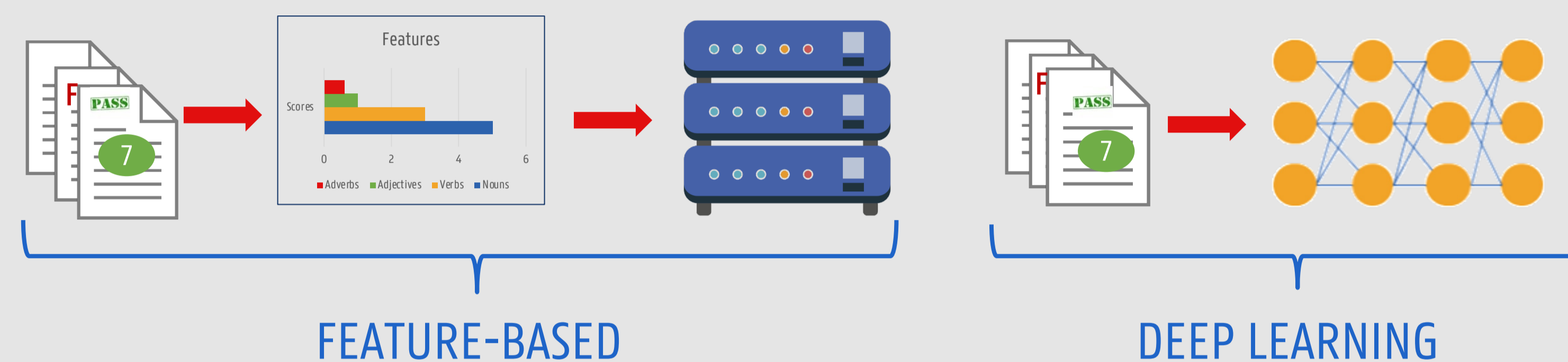
Explore creation of a Dutch AES system able to score young learners' writing skills

AUTOMATED ESSAY SCORING (AES)

- » Used for high-stakes tests → GRE & TOEFL (Richardson, 2021)
- » High correlation with human raters (Allen et al., 2016)
- » Mainly researched on English & essays (Strobl, 2019)

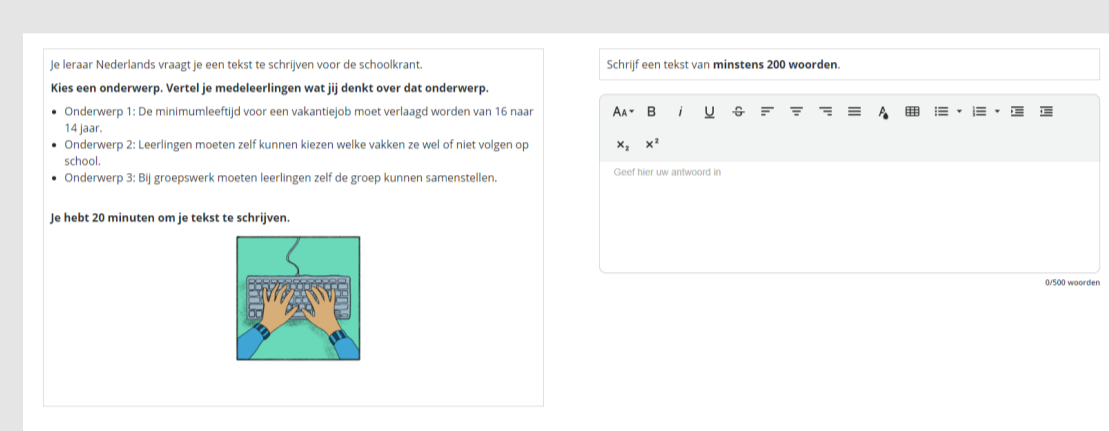


STATE OF THE ART = MACHINE LEARNING



RELIABLE DATA = PREREQUISITE TO TRAIN AES

- » Reliable → representative dataset
- » Flemish pupils 2nd year secondary education
- » Stream A and stream B
- » Different genres



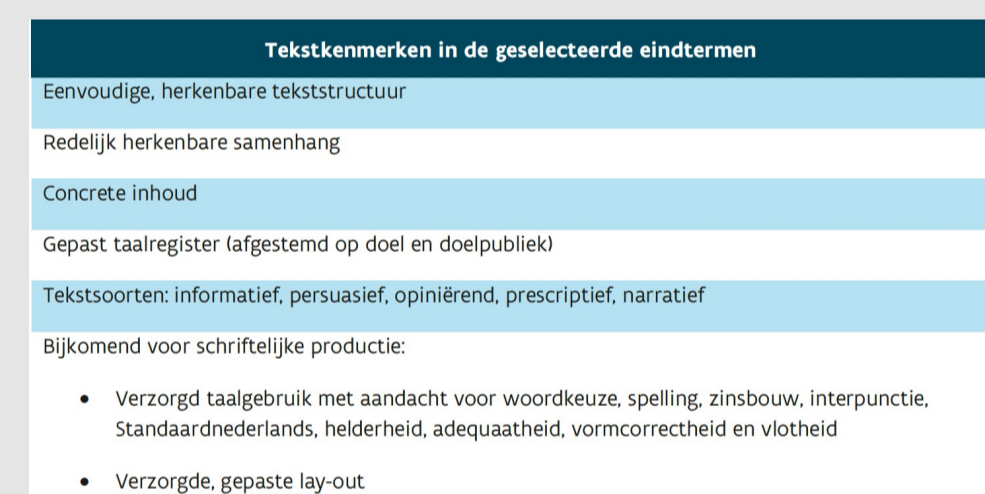
Example prompt: INF

	Informative	Opinion	Avg. word length	Min words
STREAM A	3681	3677	843	14
STREAM B	1093	1065	622	10

1,502 texts selected

HUMAN

- » Reliable → double-scoring

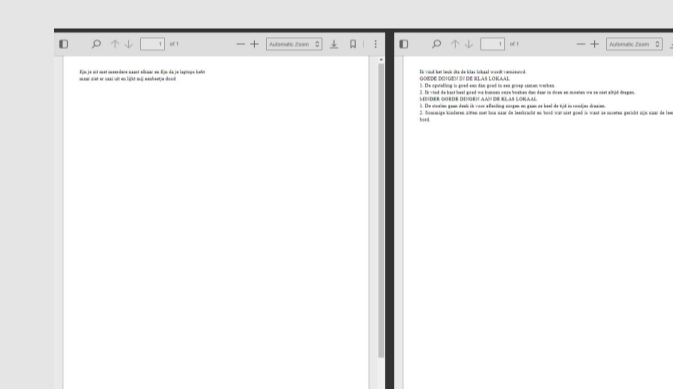


1. Analytical scoring

- » What? Compare a text on certain aspects
- » How? Designated scoring rubric + benchmarks

2. Comparative judgment

- » What? Compare 2 texts (approx. 2 hours)
- » How? Holistic scoring



- » Results expected Spring 2025

MACHINE

- » Experiments with historical data (STEP)

Assessment	# texts
BL instruction 1: E-mail (INF)	1254
BL Instruction 2: Klaslokaal (OPI)	1270
A instruction 1: Burgemeester (PERS)	684
A instruction 2: Film A (INF)	682
B instruction 1: Kippen (PRES)	564
B instruction 2: Film B (INF)	547

- » Scored with comparative judgment

FEATURE-BASED

- » 388 features extracted with T-SCAN\*

Lexical complexity Cohesion  
Grammatical complexity Fluency

- » Regression experiments

Experiment with all (388) features  
Experiment with features after LASSO regression



DEEP LEARNING

- » Fine-tuning an encoder-based LLM



» BERTje (de Vries et al., 2019)

- » Regression experiment



RESULTS

Assessment	10-fold cross validation						held-out evaluation split					
	FEATURE-BASED		DEEP LEARNING		FEATURE-BASED		DEEP LEARNING		FEATURE-BASED		DEEP LEARNING	
	All features	LASSO	LLM	LLM	All features	LASSO	LLM	All features	LASSO	LLM	All features	LASSO
BL instruction 1: E-mail (INF)	RMSE: 0.1208	QWK: 0.7403	RMSE: 0.1157	QWK: 0.7617	RMSE: 0.1177	QWK: 0.8003	RMSE: 0.1085	QWK: 0.7603	RMSE: 0.1129	QWK: 0.7584	RMSE: 0.1090	QWK: 0.8064
BL Instruction 2: Klaslokaal (OPI)	RMSE: 0.1198	QWK: 0.7856	RMSE: 0.1186	QWK: 0.7923	RMSE: 0.1091	QWK: 0.8406	RMSE: 0.1173	QWK: 0.8012	RMSE: 0.1069	QWK: 0.8154	RMSE: 0.0962	QWK: 0.8647
A instruction 1: Burgemeester (PERS)	RMSE: 0.1313	QWK: 0.5533	RMSE: 0.1243	QWK: 0.5851	RMSE: 0.1339	QWK: 0.6326	RMSE: 0.1341	QWK: 0.5463	RMSE: 0.1282	QWK: 0.5884	RMSE: 0.1174	QWK: 0.6513
A instruction 2: Film A (INF)	RMSE: 0.1069	QWK: 0.5872	RMSE: 0.1021	QWK: 0.6728	RMSE: 0.1046	QWK: 0.6628	RMSE: 0.1083	QWK: 0.6276	RMSE: 0.1033	QWK: 0.7163	RMSE: 0.1166	QWK: 0.6414
B instruction 1: Kippen (PRES)	RMSE: 0.1135	QWK: 0.7369	RMSE: 0.0993	QWK: 0.7890	RMSE: 0.1008	QWK: 0.8334	RMSE: 0.1204	QWK: 0.6598	RMSE: 0.1104	QWK: 0.7525	RMSE: 0.1090	QWK: 0.7877
B instruction 2: Film B (INF)	RMSE: 0.1241	QWK: 0.7220	RMSE: 0.1135	QWK: 0.7830	RMSE: 0.1182	QWK: 0.7897	RMSE: 0.1386	QWK: 0.6954	RMSE: 0.1319	QWK: 0.7136	RMSE: 0.1166	QWK: 0.6414

FUTURE WORK

- » How to operationalize formative aspects of writing in ML setting?
- » More ML experiments on Steunpunt data + how to scale?
- » Deep learning works well BUT focus should be on feedback and transparency (features, probing, ...)

Contact

[Orphee.DeClercq@UGent.be](mailto:Orphee.DeClercq@UGent.be)  
[Joni.Kruijsbergen@UGent.be](mailto:Joni.Kruijsbergen@UGent.be)  
[Fauve.DeBacker@UGent.be](mailto:Fauve.DeBacker@UGent.be)  
[GoedeleVandommele@UGent.be](mailto:GoedeleVandommele@UGent.be)

REFERENCES  
 • Allen, L. K., Jacovina, M. E., & McNamara, D. S. (2016). Computer-based writing instruction. In C. A. MacArthur, S. Graham, & J. Fitzgerald (Eds.), Handbook of writing research (pp. 316–329). The Guilford Press.  
 • de Vries, W., van Cranenburgh, A., Bisazza, A., Caselli, T., van Noord, G., & Nissim, M. (2019). BERTje: A Dutch BERT Model. <http://arxiv.org/abs/1912.09582>.  
 • Richardson, M., & Clesham, R. (2021). Rise of the machines? The evolving role of AI technologies in high-stakes assessment. London Review of Education, 19(1). <https://doi.org/10.14324/LRE19.1.09>  
 • Strobl, C., Ailhaud, E., Benetos, K., Devitt, A., Kruse, O., Proske, A., & Rapp, C. (2019). Digital support for academic writing: A review of technologies and pedagogies. Computers & Education, 131, 33–48.

\* <https://github.com/UUDigitalHumanitieslab/tscan>