

Architektura Transformer je velmi populární, takže může být potenciálně významné interpretovat, co ovlivňuje její výkon. Testujeme hypotézu, že model se při práci s textem spoléhá na jeho lingvistické vlastnosti. Abychom eliminovali vliv kultury na význam, používáme úlohu pracující na úrovni znaků s Transformer modelem ByT5. Dotrénujeme ByT5-small na dešifrování vět zašifrovaných pomocí textových šifer (Vigenère, Enigma). Anotujeme evaluační dataset vět pomocí publikovaných nástrojů pro NLP. Na evaluačním datasetu zkoumáme vztahy mezi lingvistickými vlastnostmi a četností chyb dotrénovaného ByT5 při dešifrování vět. Analyzujeme korelace, trénujeme ML modely na predikci četnosti chyb věty z jejich lingvistických vlastností a interpretujeme důležitost vlastností pomocí SHAP. Nacházíme malé signifikantní korelace, ale predikce četnosti chyb z vlastností selhává. Dospíváme k závěru, že identifikované vlastnosti neposkytují vhled do výkonu Transformerů.