

Vizualizácia komplexných modelov neurónových sietí s architektúrou typu Transformer je vo všeobecnosti veľmi náročná úloha, ktorá sa väčšinou rieši vizualizáciou blokov Attention a sledovaním, na ktoré slová sa tento blok zameriava. Modely Transformer ale majú veľké množstvo vrstiev, na každej vrstve majú veľké množstvo hláv Attention a každá hlava môže sledovať rôzne lingvistické znaky. My sme sa preto v tejto práci zamerali na vytvorenie programu, ktorý je určený na prehľadnejšiu vizualizáciu správania jazykového modelu GPT-2. Prišli sme so štyrmi metódami vizualizácie, ktoré skúmajú závislosti generovaných slov od prechádzajúcich slov v texte. Tieto závislosti sledujeme tak, že skúsime prvé slovo v texte vynechať alebo zameniť za podobné slovo a pozorujeme zmenu v pravdepodobnosti generovaného slova. Metódy sme vyskúšali na modele GPT-2 Medium a demonštrujeme, aké výsledky dané metódy vytvorili. Zároveň vyslovujeme hypotézy, ktoré sa pokúšajú objasniť, prečo tieto výsledky vyšli práve tak.