| Layer | Lang | Feature ID | Interpretation | Detection | | | | Fuzzing | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|
| | | | | Accuracy | F1 score | Precision | Recall | Accuracy | F1 score | Precision | Recall |
| model.layers.11.mlp | Turkish | 127127 | The highlighted tokens are predominantly proper nouns, especially names of people, places, and groups, often in the context of multicultural or multilingual event descriptions. These tokens frequently appear in various scripts and languages, and are often associated with performances, orchestras, or notable individuals. | 0.61 | 0.552 | 0.649 | 0.48 | 0.55 | 0.237 | 0.778 | 0.14 |
| model.layers.12.mlp | Turkish | 40217 | The highlighted tokens are Turkish verb roots and suffixes, often marking verb stems, tense, aspect, or person, and are frequently found at the end of words to indicate actions, states, or processes. | 0.86 | 0.837 | 1 | 0.72 | 0.9 | 0.891 | 0.976 | 0.82 |
| model.layers.12.mlp | Turkish | 45932 | The highlighted tokens are predominantly Turkish morphemes, suffixes, and noun or verb stems, often marking grammatical features such as possession, plurality, tense, or case. Many are parts of compound words, inflections, or derivational endings, and frequently appear at the end of words or as part of agglutinative constructions, reflecting the morphological richness and structure of Turkish. | 0.79 | 0.741 | 0.968 | 0.6 | 0.8 | 0.756 | 0.969 | 0.62 |
| model.layers.12.mlp | Turkish | 69606 | The highlighted tokens are predominantly Turkish verb roots and suffixes, especially those forming various tenses, voices, and moods, as well as participles and nominalizations. There is a strong focus on verb morphology, including derivational and inflectional endings that indicate person, tense, aspect, modality, and negation. These patterns reflect the agglutinative structure of Turkish, where meaning is built up through the sequential addition of suffixes to verb and noun stems. | 0.7 | 0.571 | 1 | 0.4 | 0.75 | 0.675 | 0.963 | 0.52 |
| model.layers.12.mlp | Turkish | 102905 | The important tokens are predominantly Turkish morphemes, suffixes, and stems that form or modify nouns, verbs, and adjectives, often marking tense, person, plurality, comparison, or possession. These tokens are crucial for grammatical structure and meaning in Turkish, highlighting the agglutinative nature of the language where meaning is built up through the addition of multiple suffixes to word roots. | 0.77 | 0.701 | 1 | 0.54 | 0.77 | 0.729 | 0.886 | 0.62 |
| model.layers.12.mlp | Turkish | 106294 | The highlighted tokens are Turkish morphemes, suffixes, or stems that frequently appear at the end or within words, often marking grammatical features such as tense, possession, plurality, or forming nouns and verbs. These segments are crucial for word formation and meaning in Turkish morphology. | 0.93 | 0.925 | 1 | 0.86 | 0.92 | 0.915 | 0.977 | 0.86 |
| model.layers.12.mlp | Turkish | 112037 | The highlighted tokens are predominantly Turkish verb and noun roots, often with attached suffixes, as well as some common noun and adjective stems. These roots frequently appear at the beginning or within words, and are central to the meaning and grammatical structure of the sentences, reflecting the agglutinative nature of Turkish morphology. | 0.84 | 0.818 | 0.947 | 0.72 | 0.83 | 0.809 | 0.923 | 0.72 |
| model.layers.13.mlp | Turkish | 3558 | The highlighted tokens are predominantly Turkish noun and verb phrases, often marking key semantic units such as actions, roles, attributes, or relationships. These include compound nouns, nominalizations, and verb forms with suffixes indicating tense, possession, or plurality. The tokens frequently appear at phrase or clause boundaries, and often encapsulate the main informational or functional content of the sentence. | 0.91 | 0.901 | 1 | 0.82 | 0.91 | 0.901 | 1 | 0.82 |
| model.layers.13.mlp | Turkish | 25033 | The highlighted tokens are predominantly Turkish morphemes, suffixes, and noun or verb roots, often marking plurality, possession, tense, or forming compound words. Many are parts of common noun or verb constructions, especially those denoting actions, states, or groupings, and frequently appear at word endings or as inflectional/derivational elements. | 0.68 | 0.529 | 1 | 0.36 | 0.69 | 0.597 | 0.852 | 0.46 |
| model.layers.13.mlp | Turkish | 34094 | The highlighted tokens frequently mark Turkish grammatical suffixes, numerals, time expressions (years, months, days), quantifiers, and common function words. These patterns are typical in Turkish text for expressing dates, quantities, durations, and grammatical relationships. | 0.9 | 0.896 | 0.935 | 0.86 | 0.93 | 0.929 | 0.939 | 0.92 |
| model.layers.13.mlp | Turkish | 42039 | The highlighted tokens are predominantly Turkish morphemes, suffixes, and common word stems, including noun and verb endings, possessive and case markers, and frequently used roots. These elements are essential for grammatical structure, word formation, and meaning in Turkish, often marking tense, plurality, possession, or case, and are crucial for understanding and generating morphologically rich Turkish text. | 0.94 | 0.936 | 1 | 0.88 | 0.96 | 0.959 | 0.979 | 0.94 |
| model.layers.13.mlp | Turkish | 63582 | The highlighted tokens are predominantly Turkish suffixes and endings that indicate grammatical relationships such as possession, plurality, tense, person, and case, as well as common noun and verb forms. These morphemes are essential for sentence structure and meaning in Turkish, often attached to root words to convey nuanced information. | 0.77 | 0.701 | 1 | 0.54 | 0.77 | 0.701 | 1 | 0.54 |