Table 2: Retrieval Results (● & ○ denote statistical significance at p-value < 0.01 & < 0.05 respectively)



| t=5 | | | | | t=50 | | | |
| Dataset | #QAs | Prec | SR | MAP | NDCG | Prec | SR | MAP | NDCG |
|---|---|---|---|---|---|---|---|---|---|
| *stats* | 4004 | **0.016**• | **0.076**• | **0.057**• | **0.060**• | 0.002 | 0.096 | **0.058**• | **0.071**• |
| *programmers* | 4107 | **0.020**• | **0.096**• | **0.068**• | **0.075**• | 0.002 | 0.115 | **0.069**• | **0.088**• |
| *wordpress* | 4744 | **0.019**• | **0.091**• | **0.069**• | **0.074**• | 0.002 | 0.112 | **0.070**• | **0.085**• |
| *physics* | 5044 | **0.025**• | **0.120**• | **0.088**• | **0.094**• | 0.003 | 0.148 | **0.090**• | **0.111**• |
| *mathematica* | 5084 | **0.018**• | **0.087**• | **0.067**• | **0.072**• | 0.002 | 0.116 | **0.069**• | **0.084**• |
| *unix* | 5330 | **0.023**• | **0.115**• | **0.089**• | **0.094**• | 0.003 | 0.137 | **0.091**• | **0.107**• |
| *gaming* | 6398 | **0.034**• | **0.166**• | **0.130**• | **0.137**• | 0.004 | 0.189 | **0.132**• | **0.155**• |
| *english* | 6668 | **0.024**• | **0.115**• | **0.090**• | **0.095**• | 0.003 | 0.130 | **0.092**• | **0.107**• |

Table 3: LASER-QA Results (Boldfacing and Statistical Significance indications from comparison with TopicTRLM and TBLM) over Larger Categories in CQADupStack



Figure 1: NDCG (Y-axis) v/s. $k$

creases way beyond the training neighborhood size (i.e., 15), LASER-QA is seen to deteriorate gracefully (as expected).

- LASER-QA performance peaks on rank-aware metrics such as MAP and NDCG (even at $t = 50$), indicating it's high effectiveness in producing relevant results at the top.

These observations underline the effectiveness of LASER-QA as a CQA retrieval method. It may be noted that LASER-QA uses compact representations ($d < 2000$), as compared to vocabulary space representations that are typically $\geq 5000$.

**Trends at High t:** The performance trends at high values of $t$ are explained by the usage of the local neighborhood (of size $k$) in LASER-QA latent space learning. Going down the result list much beyond $k$ reveals expected, but graceful, decline in accuracy. For automated processing scenarios that necessitate large $t$, a correspondingly high $k$ may be used in learning. It is notable that LASER-
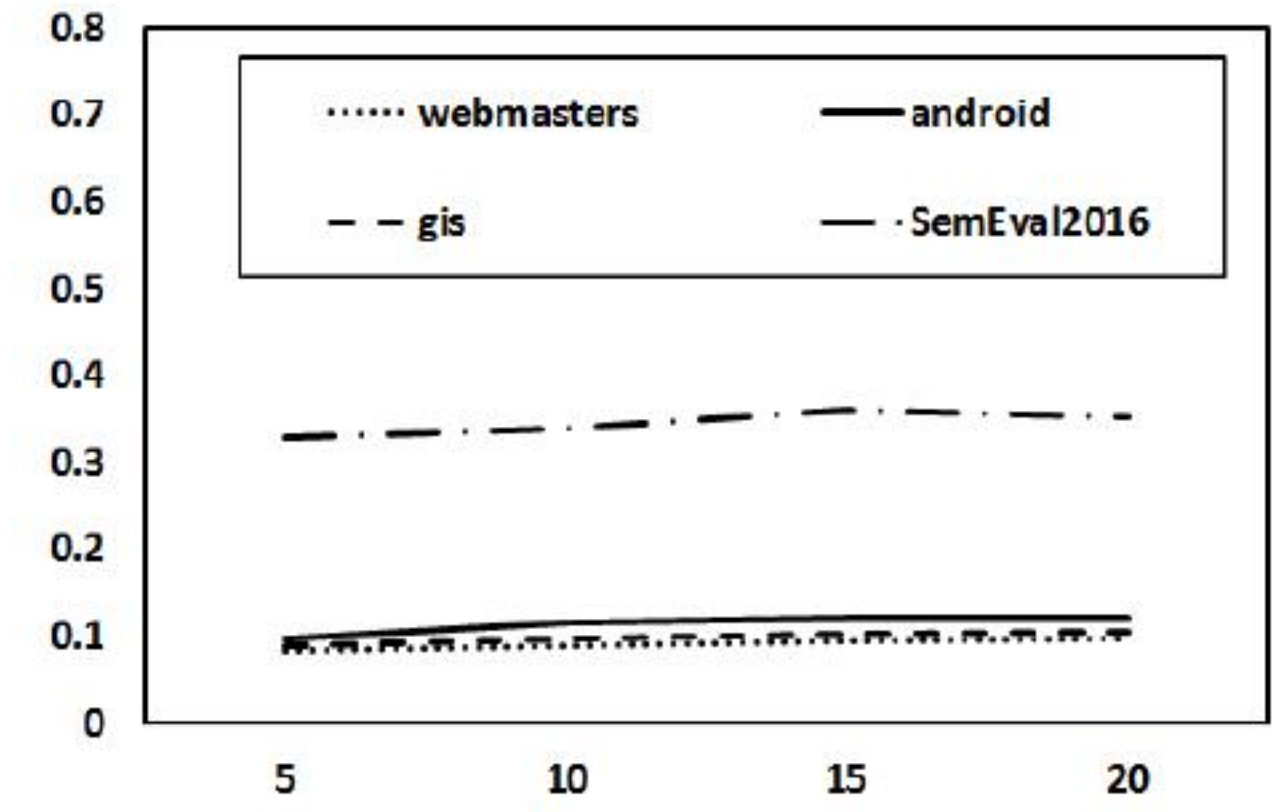
QA's focus on local neighborhood manifests as enhanced accuracy at the top of the result set.

**LASER-QA Analysis on Larger CQADup-Stack Datasets:** Owing to scalability issues of AENN that disallows a full evaluation over larger categories in CQADupStack, we present LASER-QA results over them in Table 3 to illustrate the consistency in trends. Boldfacing and statistical significance have the same semantics as earlier, with the comparison performed against only TopicTRLM and TBLM.

### 5.3 LASER-QA Parameter Analysis

We now analyze the NDCG trends (NDCG being the most popular IR metric) across LASER-QA parameters, i.e., $k$, $\alpha$ and $d$, varying each one separately keeping the default choice for others.

- **Varying $k$:** Figure 1 plots NDCG against values of $k$ from $\{5, 10, 15, 20\}$. As may be seen, the accuracy is seen to improve with increasing $k$ in the lower ranges, while sat-