



**Zadatak** Koristeći Apache Lucene i Tika biblioteke indeksirati kolekciju dokumenata različitih formata i odraditi pretraživanje nad kreiranim indeksom.

1. Kreirati Java aplikaciju u sa nazivom **<broj\_indeksa>\_<ime>\_<prezime>\_lab1** u Eclipse okruženju i referencirati potrebne JAR fajlove iz Apache Lucene biblioteke i iz Tika biblioteke.
2. Pronaći najmanje 3 dokumenta različitih formata (doc, docx, pdf, html, ppt, pptx,...) veličine od 30KB do 1MB. Fajlove u nekim od ovih formata je moguće pronaći npr. na sajtu Project Gutenberg ([www.gutenberg.org](http://www.gutenberg.org)) koji sadrži preko 50000 besplatnih knjiga **Obavezno je da svi studenti imaju različite kolekcije fajlova sa kojima rade.**
3. Indeksirani dokument zbog ograničenja u veličini ne mora da bude kompletna knjiga. Možete indeksirati poglavlje knjige, PPT prezentaciju, neki tekst u HTML ili DOC formatu...
4. Kreirati indeks nad tim fajlovima tako da svaki fajl predstavlja poseban dokument u indeksu i sačuvati ga negde na fajl sistemu. Indeks treba da sadrži najmanje sledeća polja: sadržaj fajla, naziv fajla, kompletnu putanju do fajla na fajl sistemu i veličinu fajla u bajtovima.
5. Kreirati bar jedan logički upit od najmanje 3 termina gde će se primeniti sve 3 logičke operacije (**NOT**, **AND** i **OR**) i izvršiti ga nad jednim i nad drugim indeksom. **Identični upit kreirati na 2 načina** – direktno kreirati objektni model bez korišćenja parsera i isti objektni model upita kreirati parsiranjem tekstualnog upita.
6. Kreirati još jedan upit. Tip upita odredite na osnovu **poslednje cifre svog broja indeksa**:
  - 0 ili 5 – **TermRangeQuery**,
  - 1 ili 6 – **PointRangeQuery**,
  - 2 ili 7 – **PrefixQuery**,
  - 3 ili 8 – **WildcardQuery**,
  - 4 ili 9 – **PhraseQuery**.

Kreirani upit izvršiti nad jednim i nad drugim indeksom. **Identični upit kreirati na 2 načina** – direktno kreirati objektni model bez korišćenja parsera i isti objektni model upita kreirati parsiranjem tekstualnog upita.