

Rangirajući pretraživač dokumenata - Infinity

Marko Budiselić

30.11.2015.

1 Predprocesiranje

2 Algoritmi

2.1 Bag of words

2.2 Vector space

2.3 Binary independence

3 Upute za pokretanje

<https://infinity.buda.link>

TODO: source setup.py

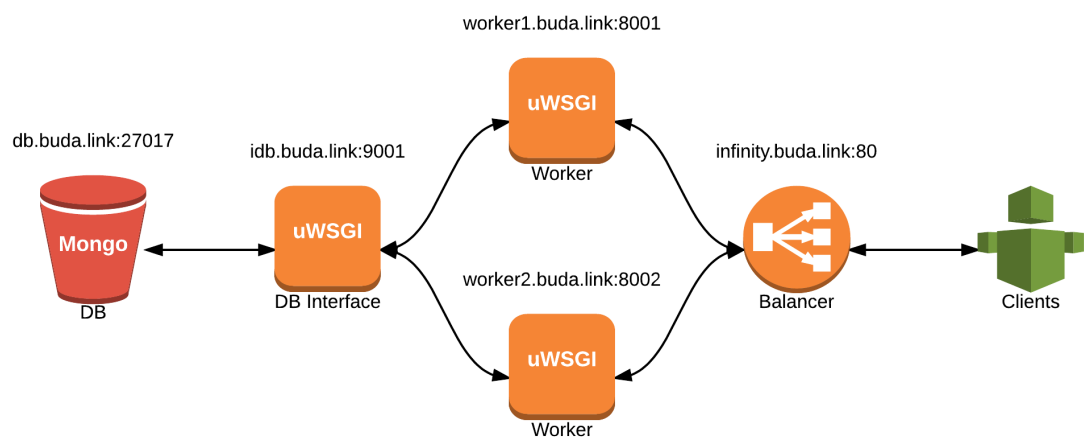
4 Produkcijska okolina

Svi upiti klijenata dolaze na Nginx load balancer koji ima izrazito veliku propusnost. Nakon toga load balancer prosljeđuje upit na worker instance koje svaka za sebe imaju cijeli dataset i znaju vratiti odgovarajući rezultat. Dataset worker instance preuzimaju od db interface instance, a ne direktno iz baze. Trenutno je u deploymentu samo jedna instanca sučelja prema bazi, ali u produkciji tu može biti opet load balancer i više instanci interface-a prema bazi podataka. Baza podataka je MongoDB, također samo jedna instanca, no u praksi tu može doći mongo replica set ili mongo shard cluster.

Kao što se vidi na grafu 1, sve instance imaju simbolička imena (FQDN). To je također izrazito bitno jer se time postiže transparentnost pristupa i migracijska transparentnost. U konkretnoj implementaciji sva imena su definirana u `/etc/hosts` file-u na deploy stroju, no u produkciji će imena biti definirana na redundantnim DNS serverima.

P.S.

infinity > gugol



Slika 1: Produkcijaska okolina