

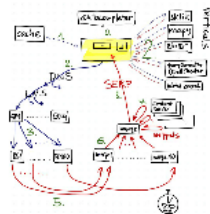
Roman Rožník seznam fulltext revealed



query processor



the Search



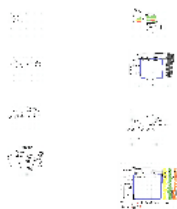
crawler



indexing



machine learning



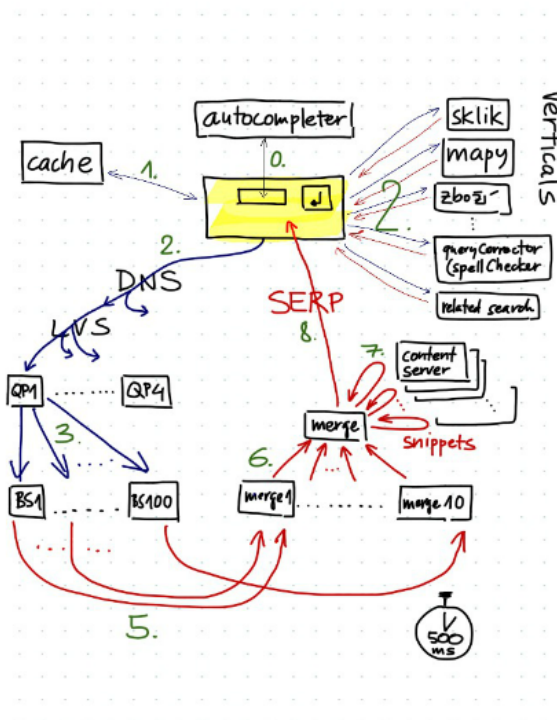
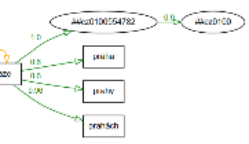
reverse index



конец

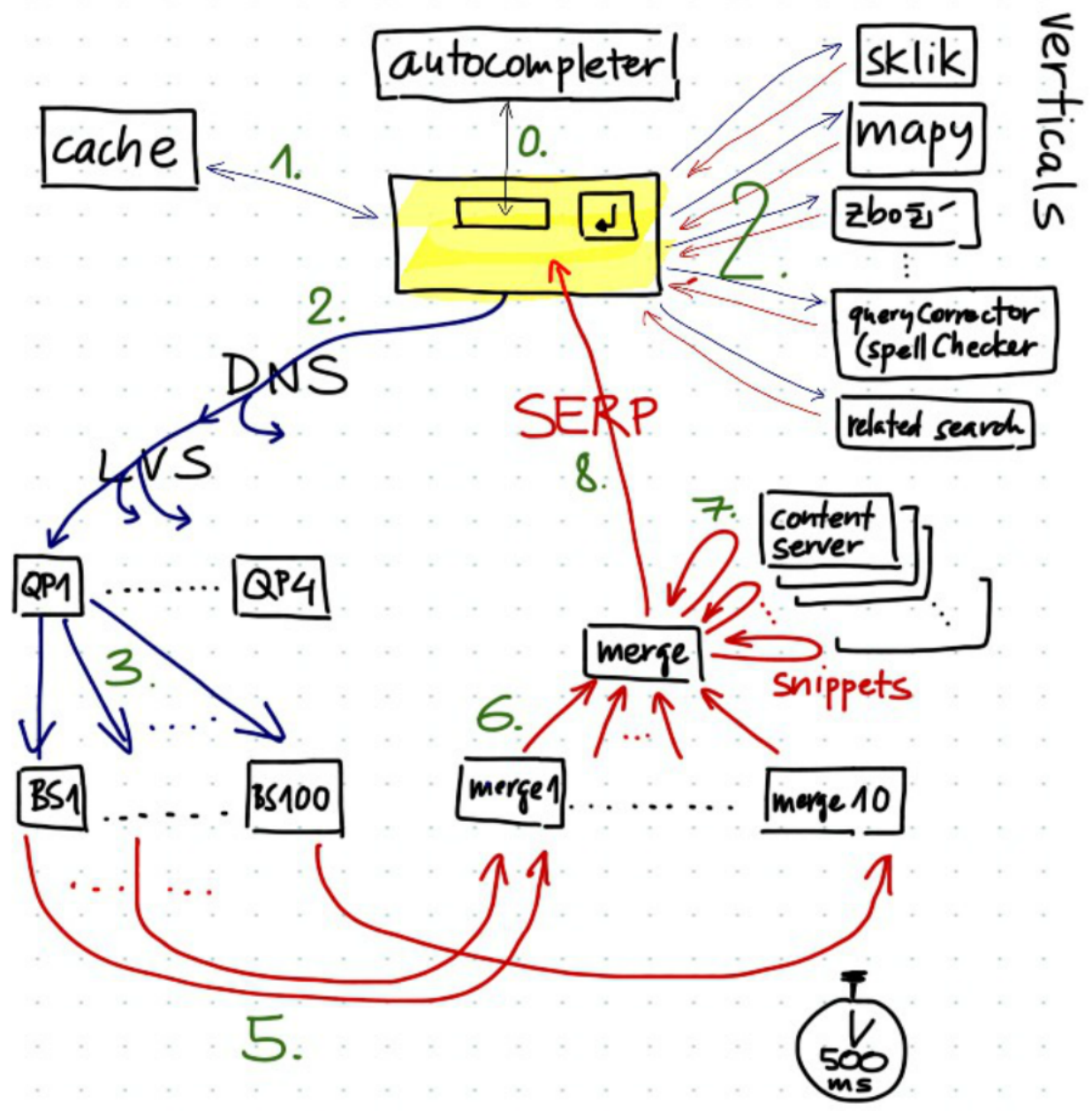
roznik@gmail.com
roman.roznik@firma.seznam.cz
kubicek@dog.cz
<https://www.youtube.com/user/RomanRoznik>
<http://bybit.fulltext.seznam.cz/seznam/>
seznam.cz

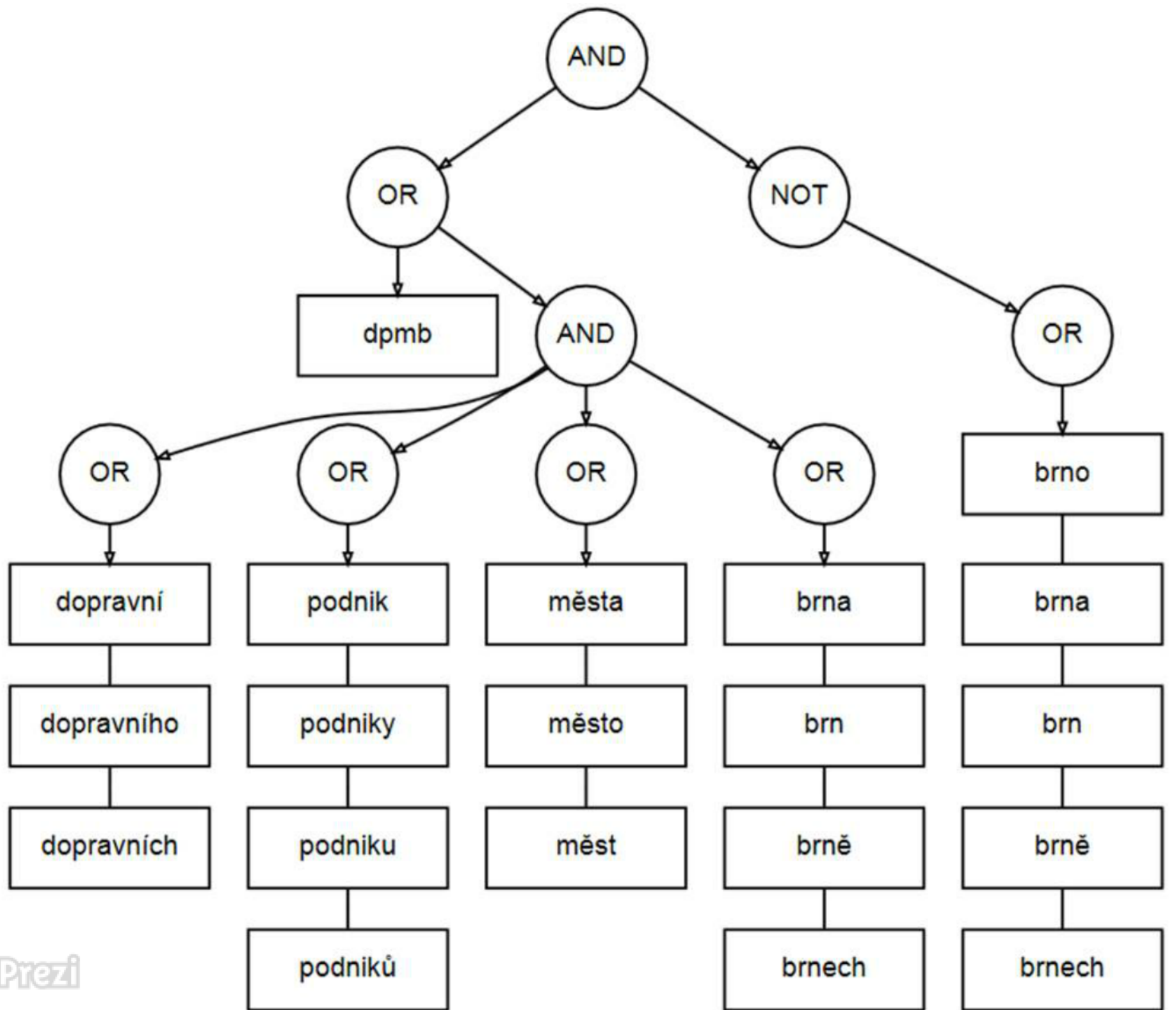
the Search

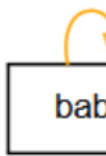
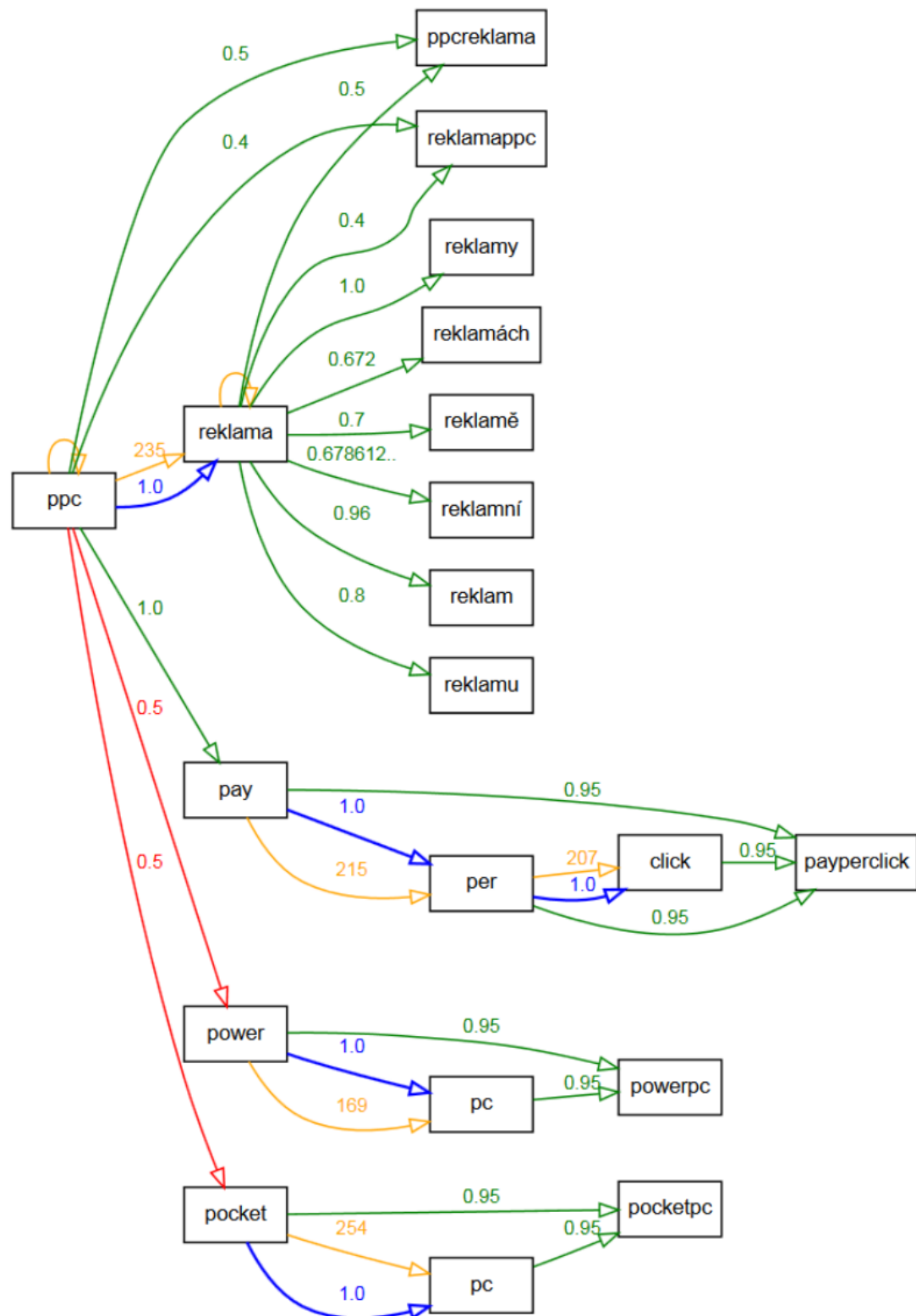


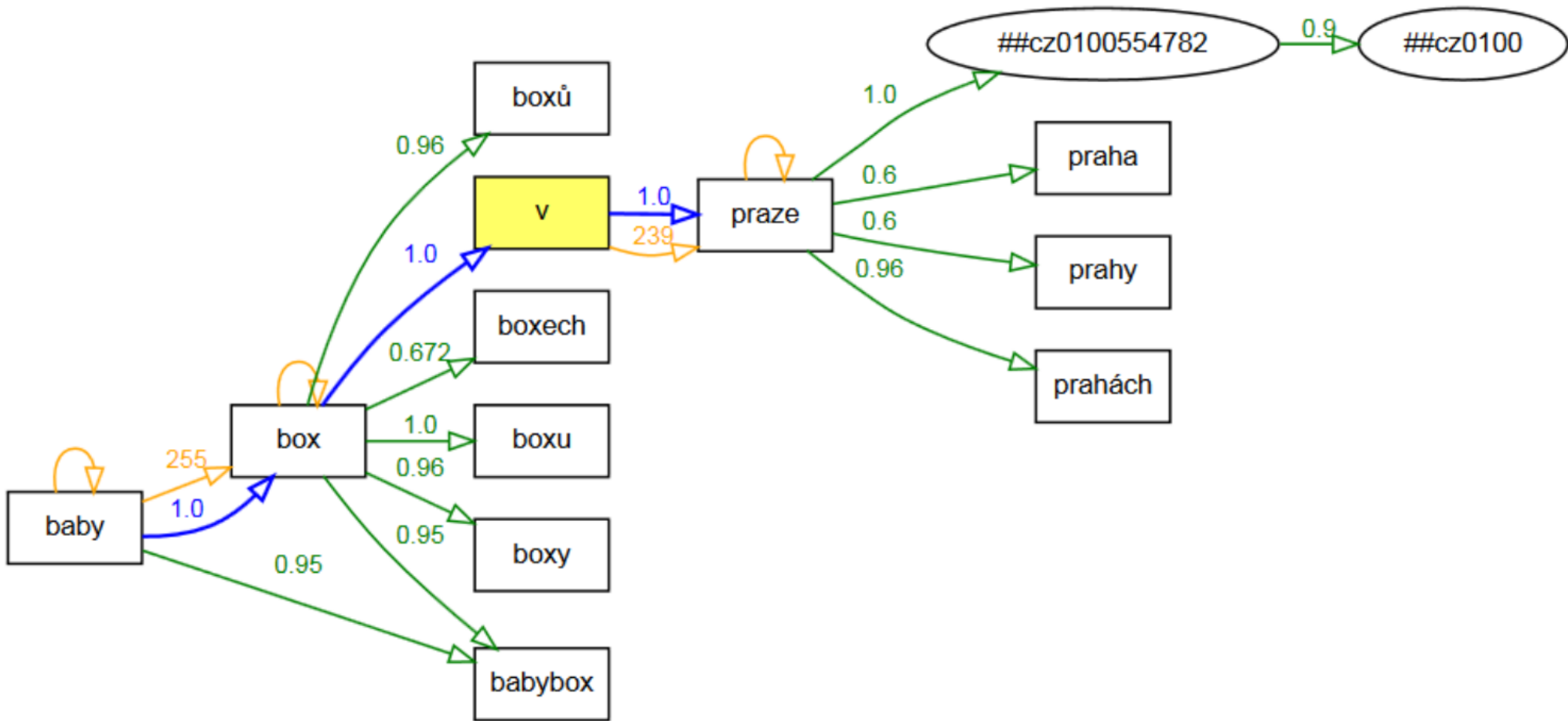
arning

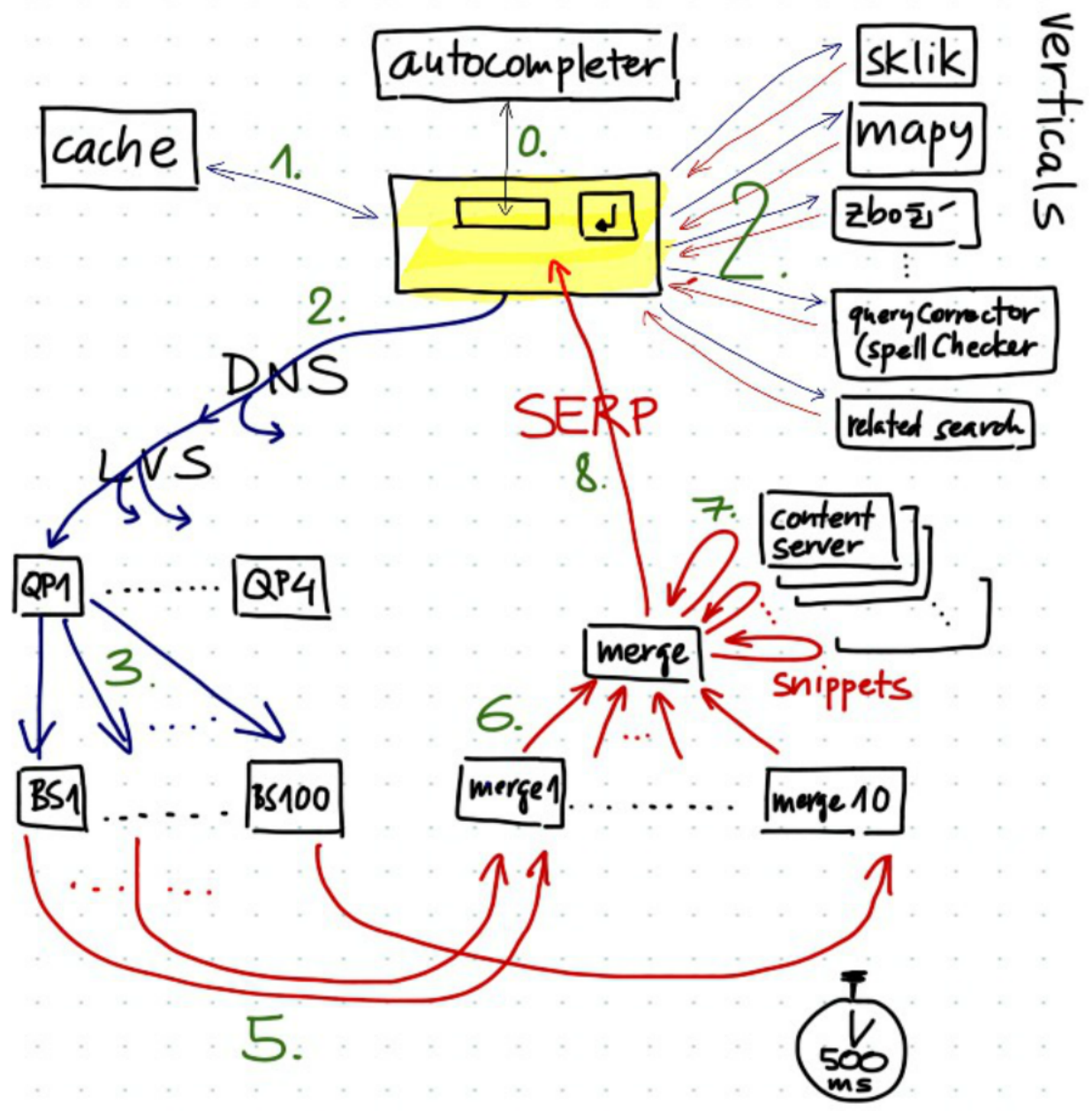




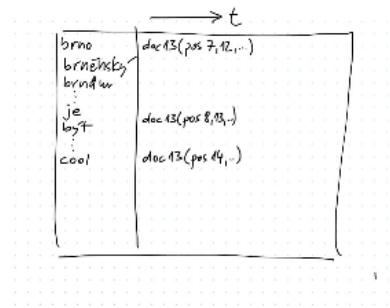




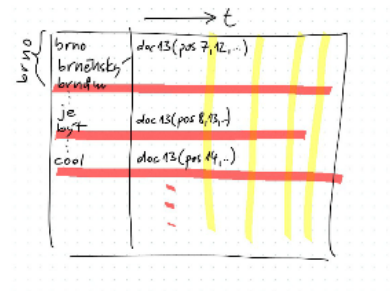




reverse index



RAM, disk, ...



→ t

brno	doc 13 (pos 7, 12, ...)
brněňsky	
brněn	
...	
je	doc 13 (pos 8, 13, ...)
byť	
...	
cool	doc 13 (pos 14, ...)

1

→ t

brno

brno

brněnský

brněn

⋮

je

byť

⋮

cool

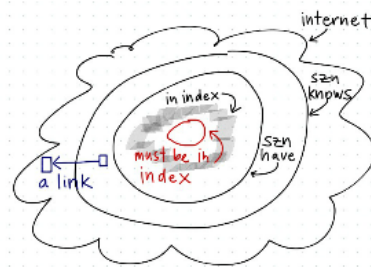
doc 13 (pos 7, 12, ...)

doc 13 (pos 8, 13, ...)

doc 13 (pos 14, ...)

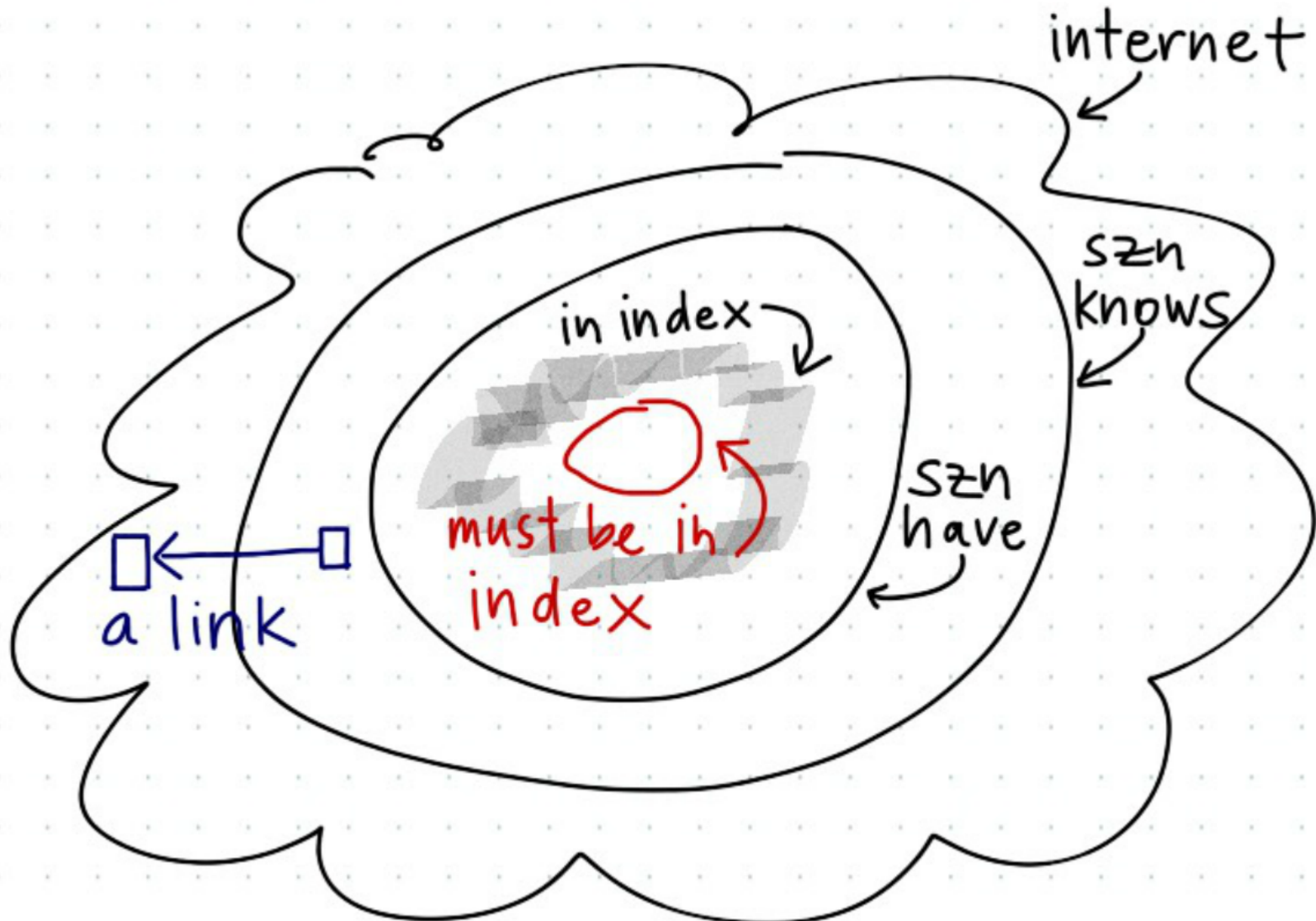
⋮

crawler



www.fi.muni.cz/usr/pelikan

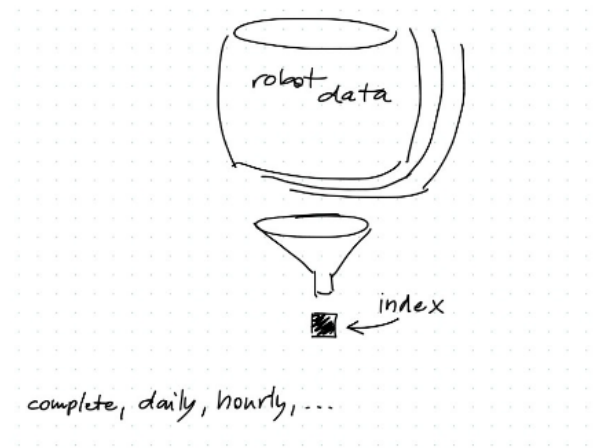
1.1.1970	500
2.3.1998	200
2.9.1998	304
2.9.1999	304
2.9.2002	304
2.9.2008	304
≈ 2020	planned

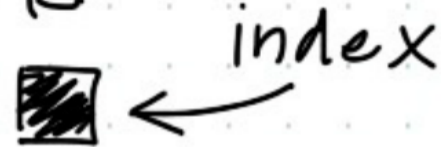
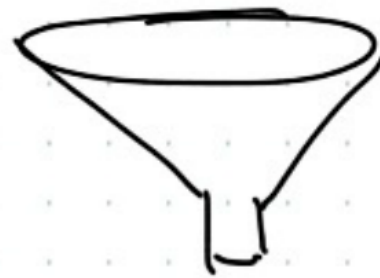


www.fi.muni.cz/usr/pelikan

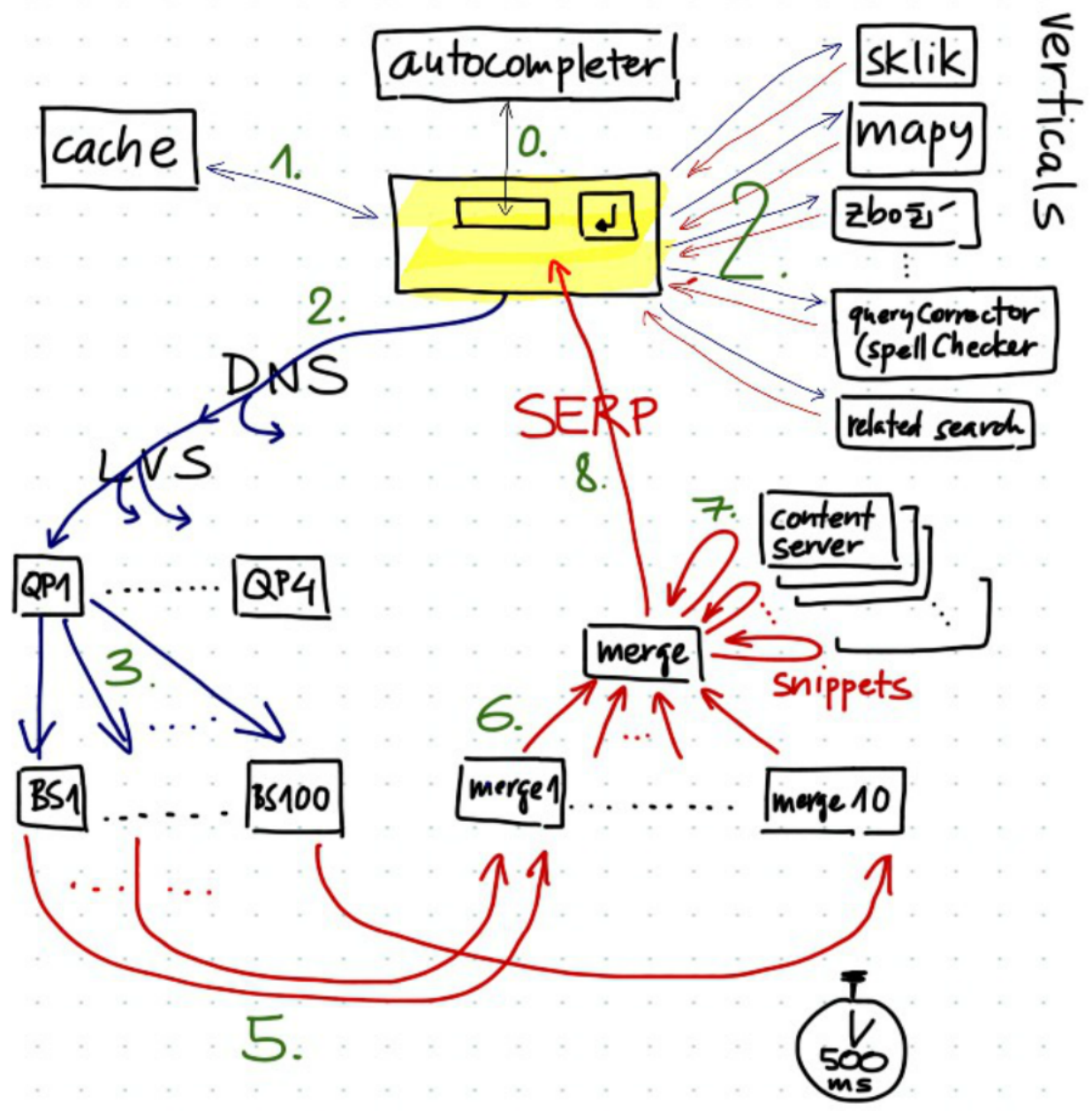
1.1.1970	500
2.3.1998	200
2.9.1998	304
2.9.1999	304
2.9.2002	304
2.9.2008	304
≈ 2020	planned

indexing

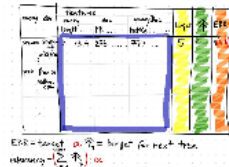
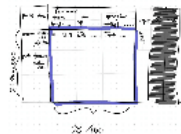
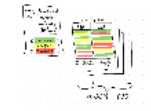




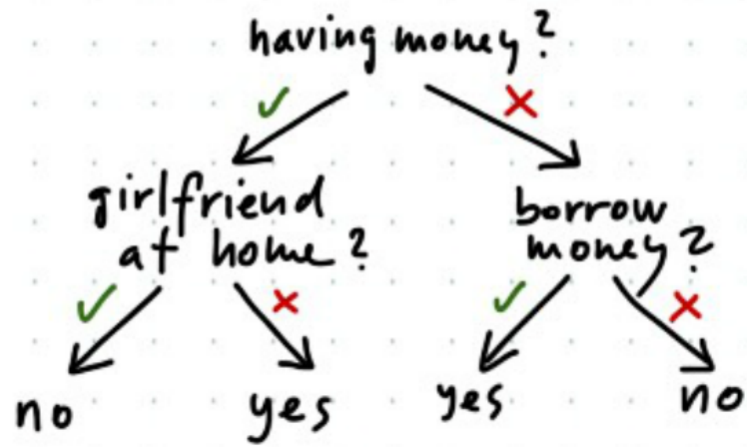
complete, daily, hourly, ...



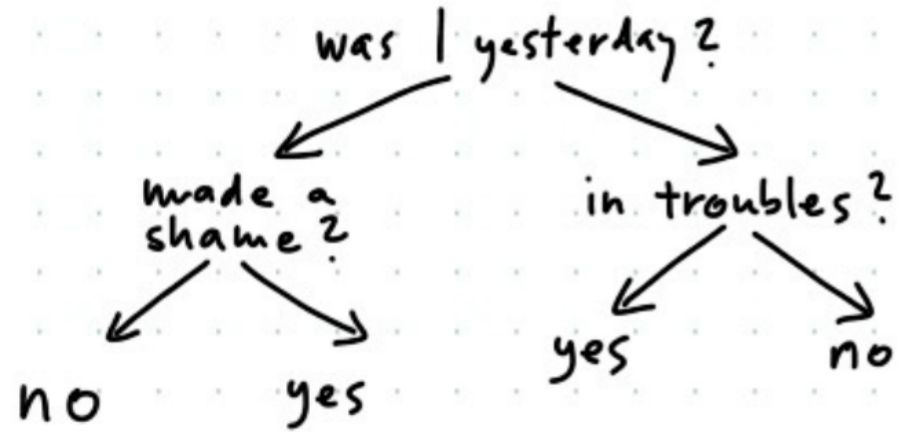
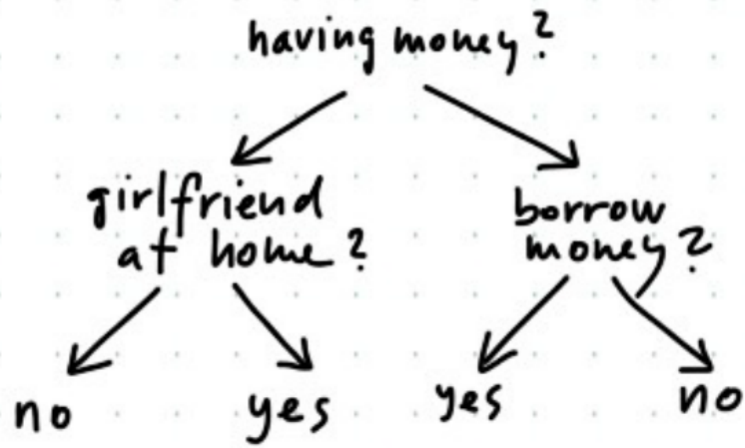
machine learning



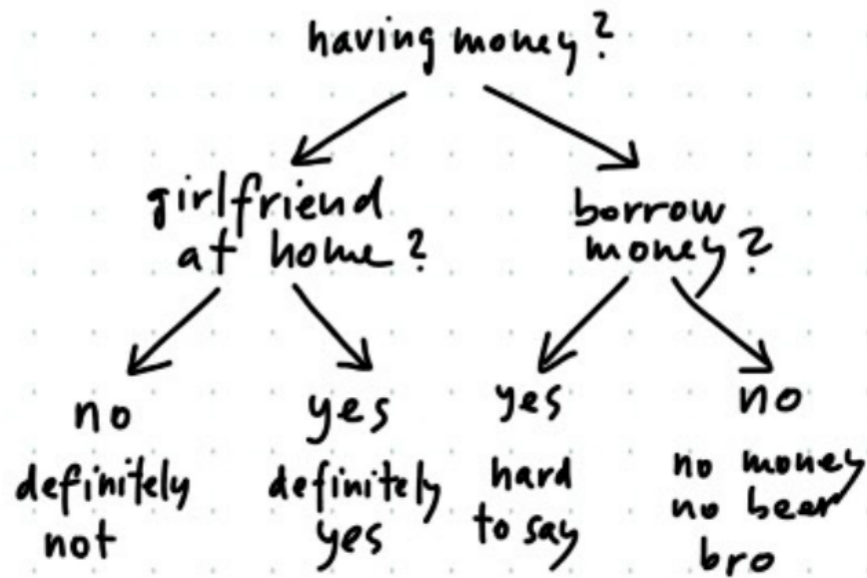
beer today?



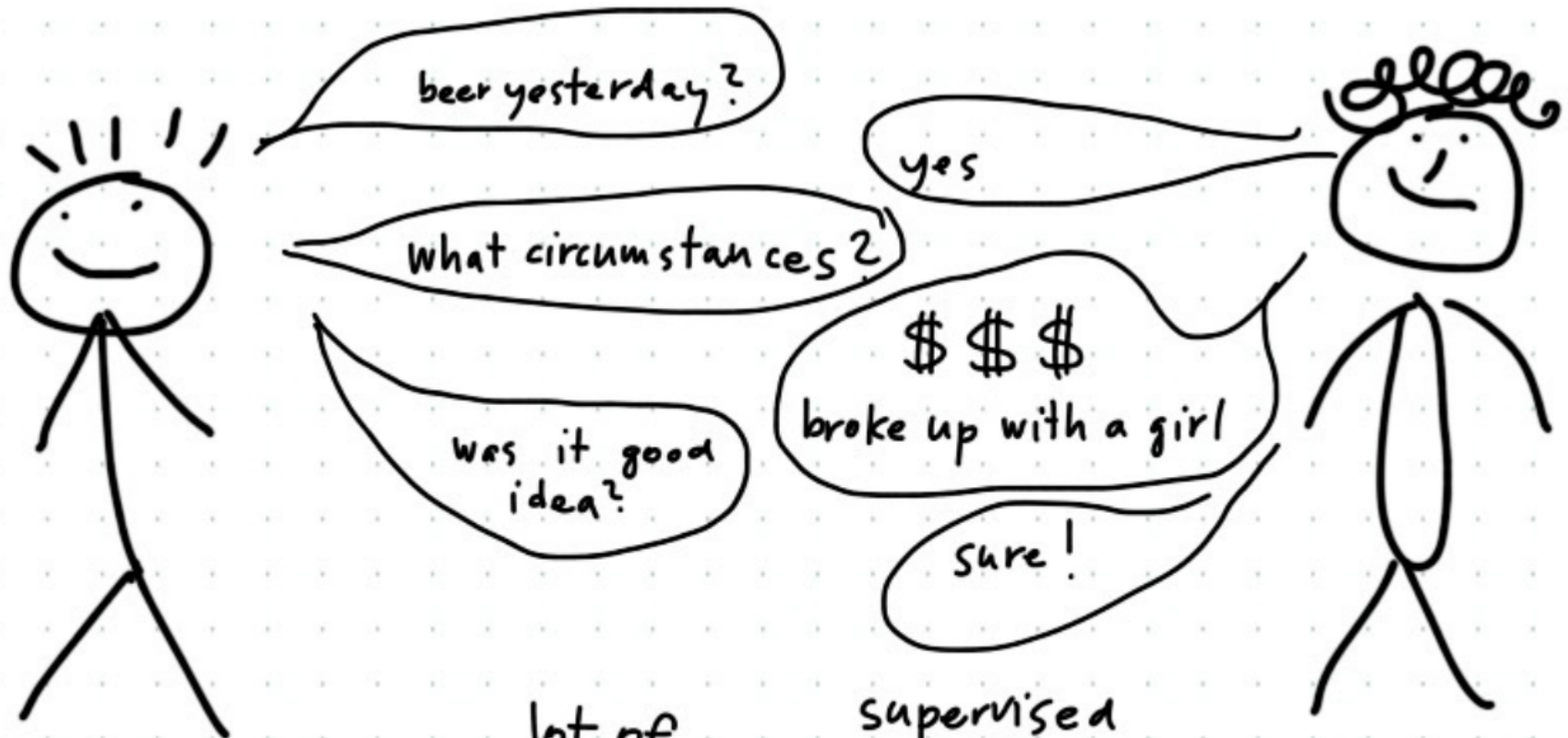
beer today?



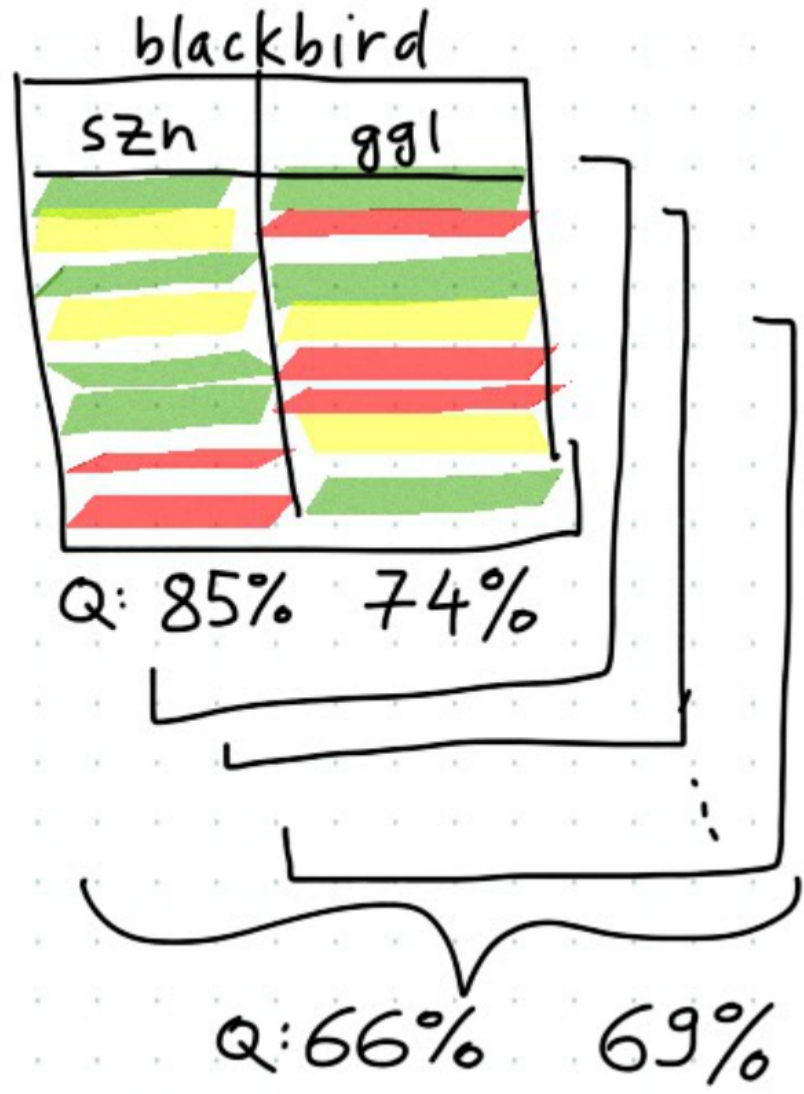
beer today?



data, data, other data,
yet another data,
even more data



lot of data → supervised machine learning

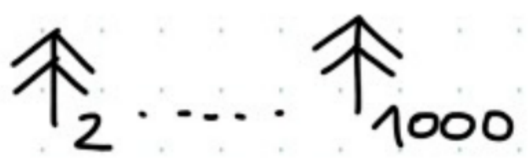
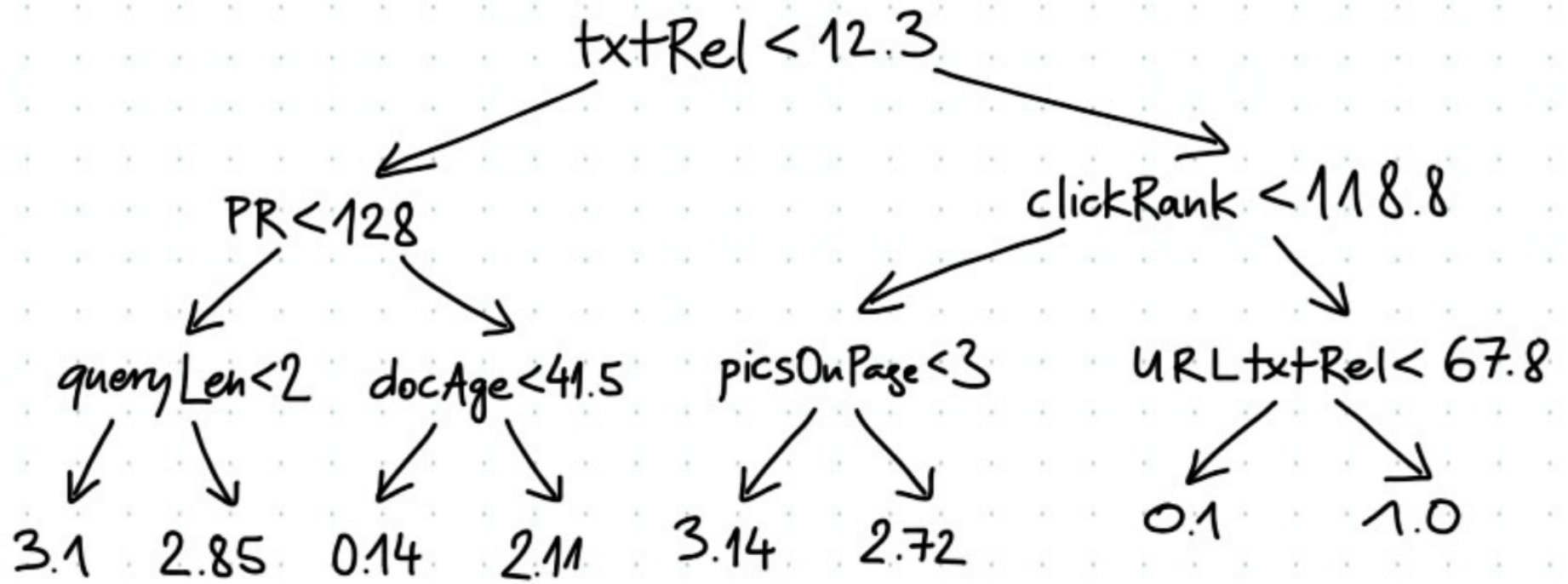


query doc	features			target	ERR
	query length...	doc PR ...	query,Doc txtRel		
seznam .cz seznam .cz wiki/sr	1	13.4	256	5	
porn freer.cz redtube.com				4	
...				...	
...					
...					

~1000000

~100

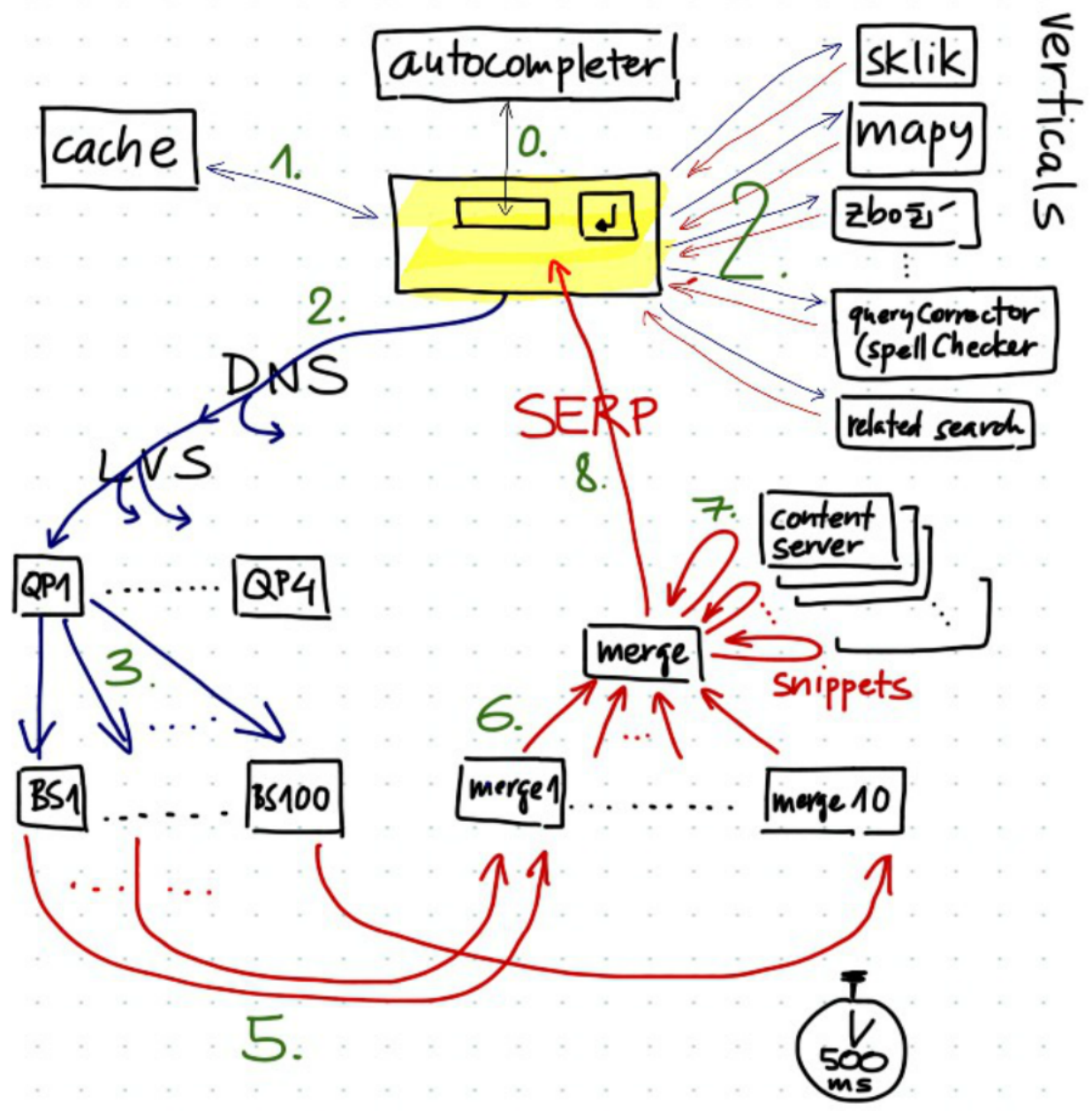
~100



query doc	features			target	↑	ERR
	query	doc	queryDoc			
	length ...	PR ...	txtRel ...			
seznam seznam .cz wiki/su	1 13.4	256	777 ...	5	4	1(4.6)
porn freer.cz redtube. com						
...						
...						

$ERR = target - \alpha \cdot \hat{\uparrow}_1 = target \text{ for next tree}$

$$relevancy = \left(\sum_{i=1}^{\infty} \hat{\uparrow}_i \right) \cdot \alpha$$



КОНЕЦ

roznik@gmail.com

roman.roznik@firma.seznam.cz

fulltext.sblog.cz

<https://www.youtube.com/user/DusanJanovsky>

[http://cyber.felk.cvut.cz/research/theses/
papers/471.pdf](http://cyber.felk.cvut.cz/research/theses/papers/471.pdf)