



Wunderlist

Nagy adat felhők között

Molnár Dániel
Microsoft



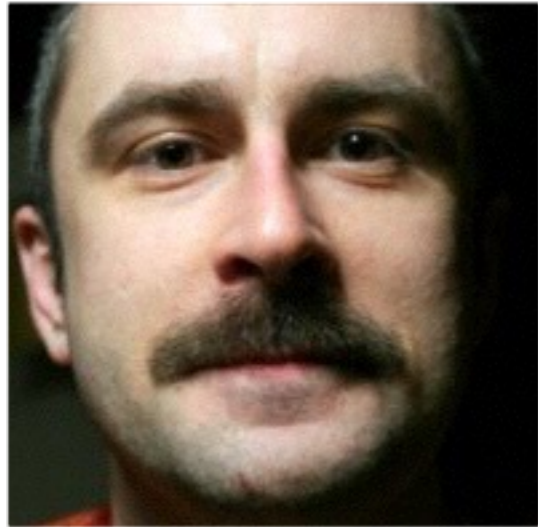
Az Úr hangja



I am the Metatron.

Az adat csapat

Univerzális emberek szűk csapata



Molnár Dániel
Infrastruktúra
BI + ML



Faludi Bence
Infrastruktúra
BI + ML



Bónis Balázs
Infrastruktúra
BI + ML



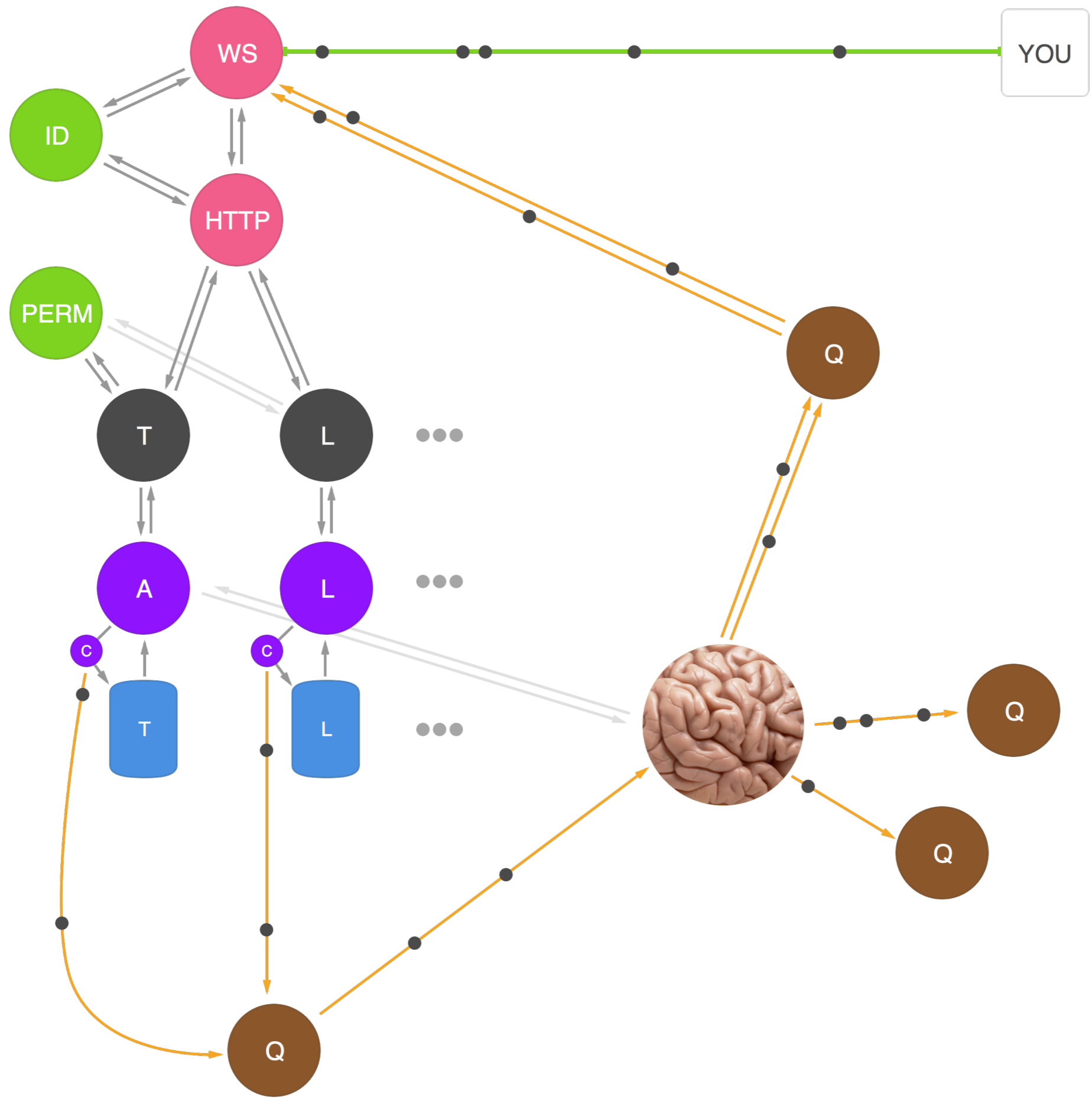
Arno Franz
BI + PM

Min dolgozunk?

A 6Wunderkinder Berlinben készíti a Wunderlistet

- ▶ Feladatkezelő iPhone, iPad, Mac, Android, Windows, Kindle Fire és web platformokon
- ▶ 21 millió felhasználó 6 év alatt 67 emberrel
- ▶ Monolitikus Rails alkalmazásból többnyelvű microservice-ek - Scala, Clojure, Go AWS-en
- ▶ 1 éve a Microsoft OneNote csapat része









Honest Status Page

@honest_update



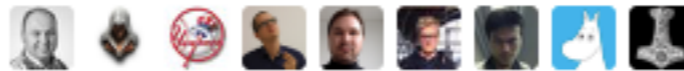
 Követés

We replaced our monolith with micro services so that every outage could be more like a murder mystery.

 Fordítás megtekintése

RETWEET
2 652

KEDVELÉS
1 929



16:10 - 2015. okt. 7.



A méret a lényeg

Produkciós adatok PostgreSQL-ben

- ▶ AWS
- ▶ ~33 DB
- ▶ ~120 konkurrens kapcsolat/db
- ▶ általában 2-3 tábla/db
- ▶ tasks táblában 1 milliárd rekord
- ▶ naponta 35 Gb inkrementális

Logok

- ▶ 125M kliens esemény naponta
- ▶ 375GB szűrt szerveresemény naponta

Egyebek

- ▶ külső források (app store, fizetés)
- ▶ KPI-k, aggregátumok, üzleti logika:
200+ lekérdezés
- ▶ önkiszolgáló adat mindenkinek

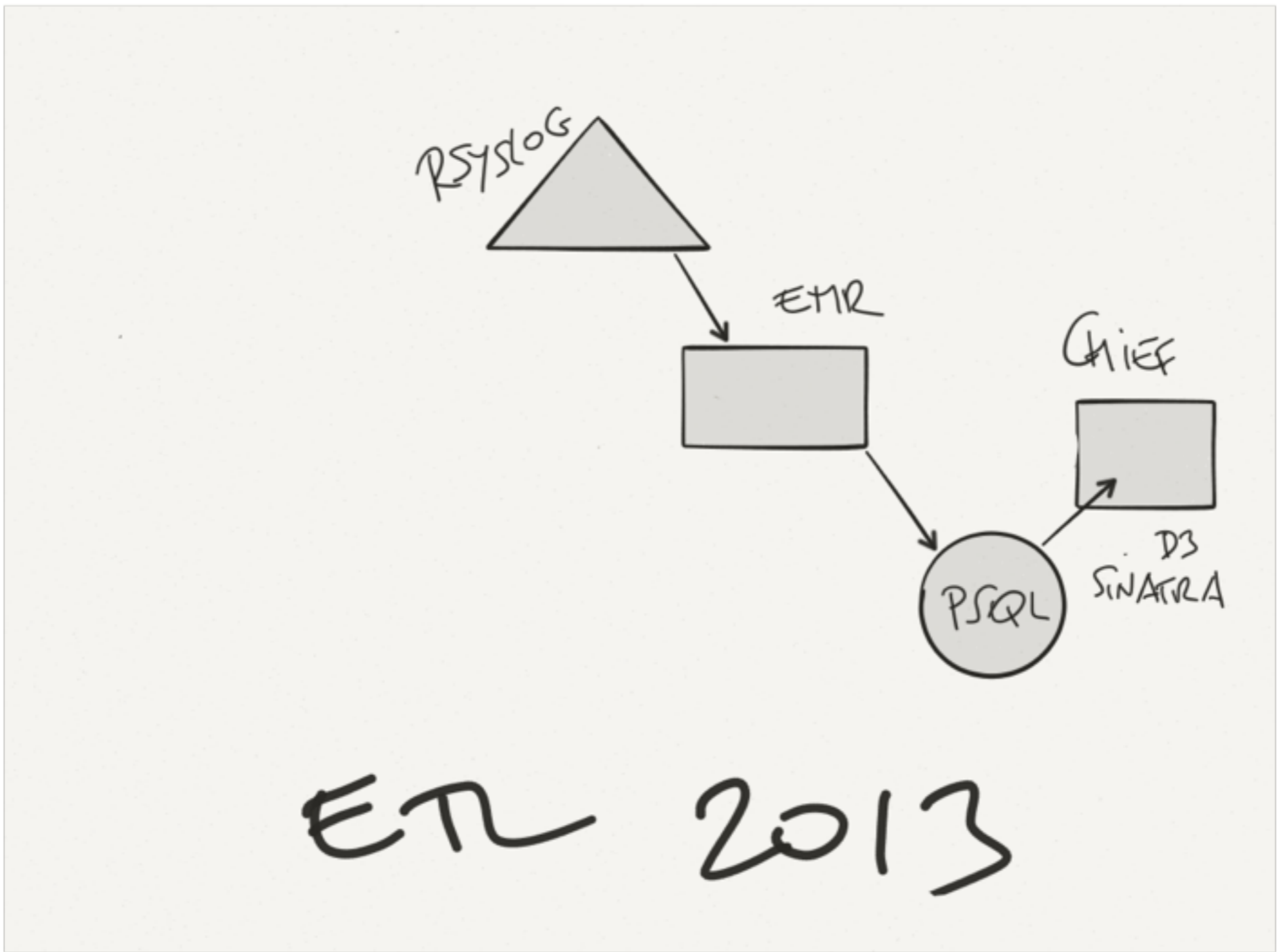


Unix + cronjob + make + SQL

Használd uncsi technológiát

Az ártatlanság kora

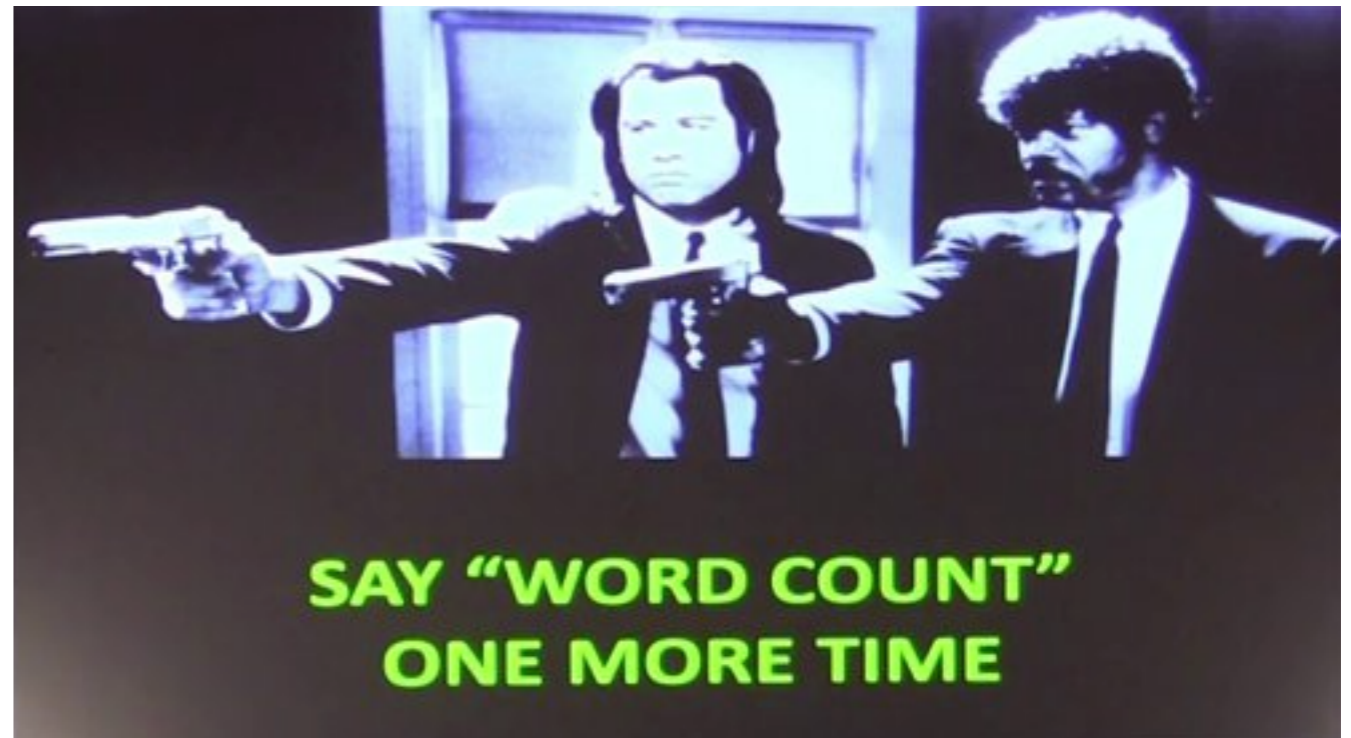
Ó, szívem, térképcsökkentés



Logok

Mondd még egyszer, hogy "Google Analytics"

- ▶ Mondj nemet! GA (nincs forrás, nincs attribúció, mintavétel, X% tévedés), Kinesis, Snowplow
- ▶ Eszközök: Railslog, Noxy, Mr Beaver (EMR Scalá-ban)
 - > Jr Beaver
- ▶ Követő node.js-ben > SNS > SQS
 - > Hamustro
- ▶ Az elosztott rendszerek logolása nagy buli (Monitorama PDX 2014 James Mickens)



ETL

Hol a WC pumpa? Mit tudom én.

- ▶ Mondj nemet! Amazon Data Flow, Oozie, Luigi
- ▶ Napi batch cronjob + make + 240 ETL SQLs
- ▶ 41 forrás (események, produkciós DB-k, App Annie, Mailchimp, fizetés, Maxmind)
- ▶ Változókat és logikát injektálunk SQL-be ERB-vel
- ▶ Bash wrapper méri a futásidőt



Night-shift

I Was Made For Lovin' You, Baby!

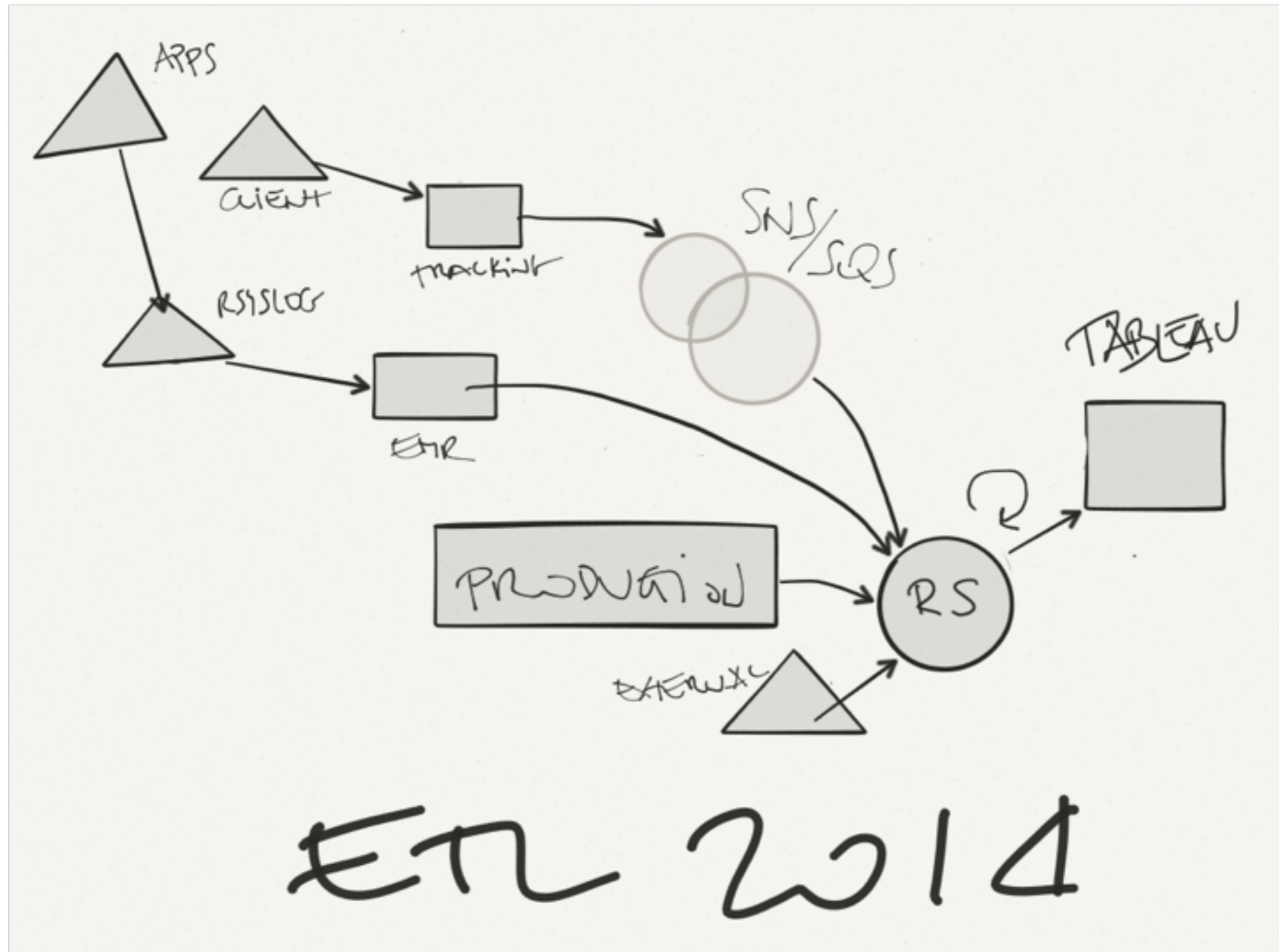
- ▶ Jeff Hammerbacher csak viccelt
- ▶ géppel olvasható dokumentáció
- ▶ csak kóddal módosítható, változások a Git-ben
- ▶ függőségek és újrapróbálkozások kezelése
- ▶ könnyen, helyi gépen tesztelhető
- ▶ párhuzamos futás
- ▶ cron időzít, bash ragaszt
- ▶ változók és logika SQL-be Ruby ERB-vel
- ▶ tracking shell méri az időt, gyűjti a hibát
- ▶ Flask interface monitorozásra
- ▶ nyílt forrású

<https://github.com/wunderlist/night-shift>



Kapcsold az ötödik sebességet

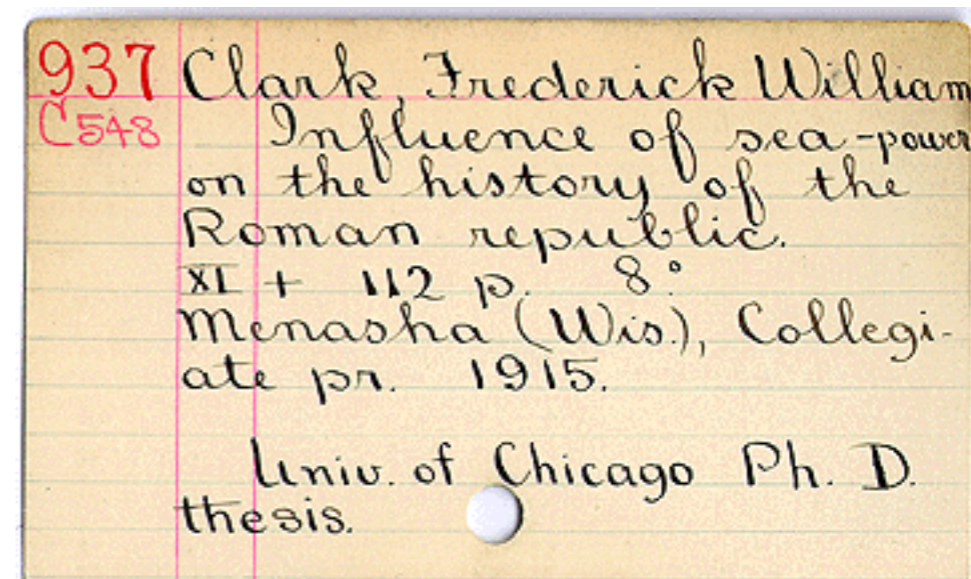
Már túl sok a nyíl



Adattárház

Dobd el a (töltött) krumplit, ha forró

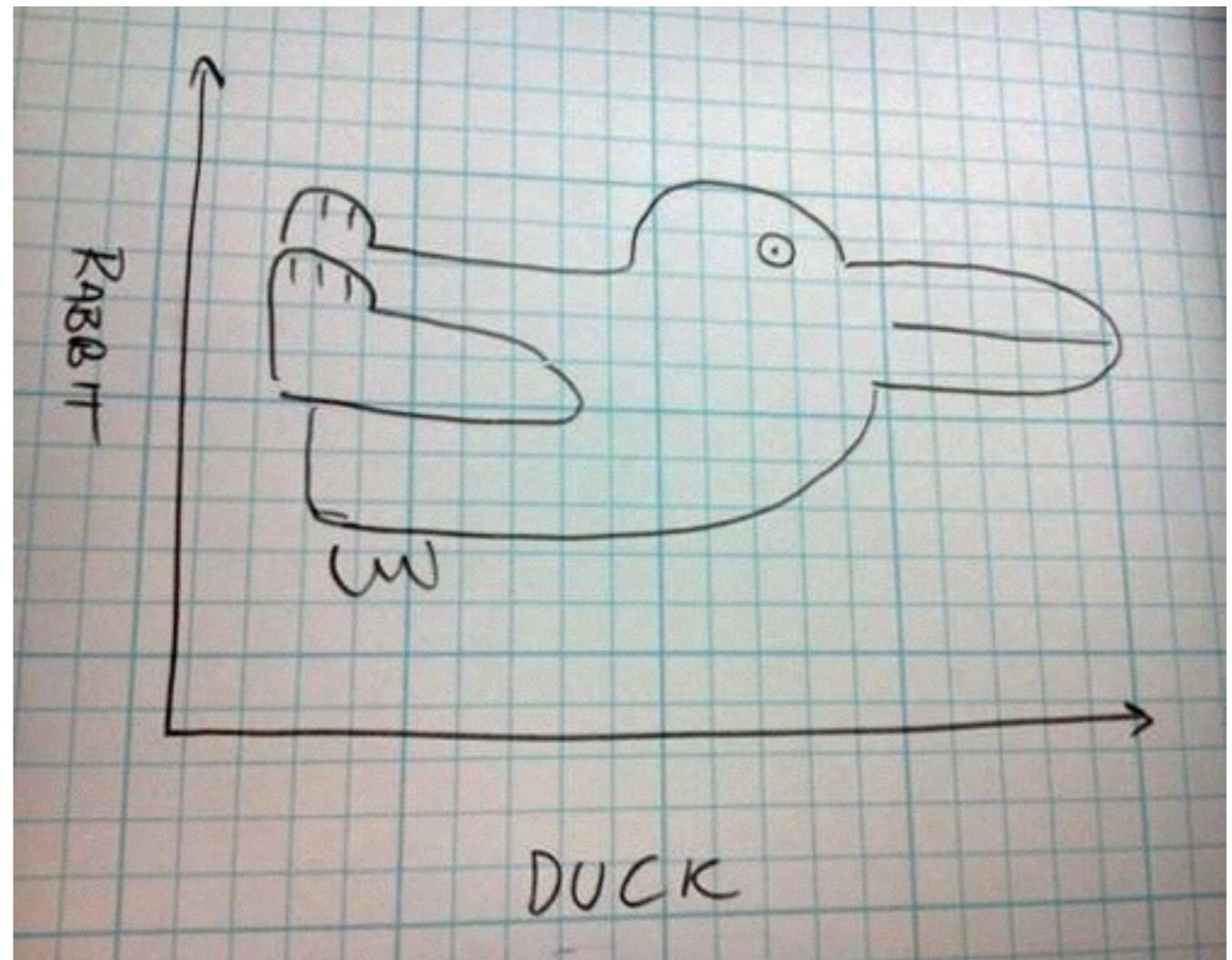
- ▶ Mondj nemet! Hadoop, Hive, Impala
- ▶ Eszközök: (PSQL) Redshift > Azure DWH
- ▶ Barebone DW, JSON, window functions
+ krek (szupergyors, olcsó)
 - GUI, támogatás, újraindítottad már?
- ▶ Ne join-olj, filterezz - minek a csillagséma
- ▶ 32 kis SSD-ről
6 kis HDD-re (cold) + 22 kis SSD-re (hot)
- ▶ 10 TB 280 táblában
- ▶ 'real' séma



Más SQL megközelítés

Attól függ, honnan nézzük

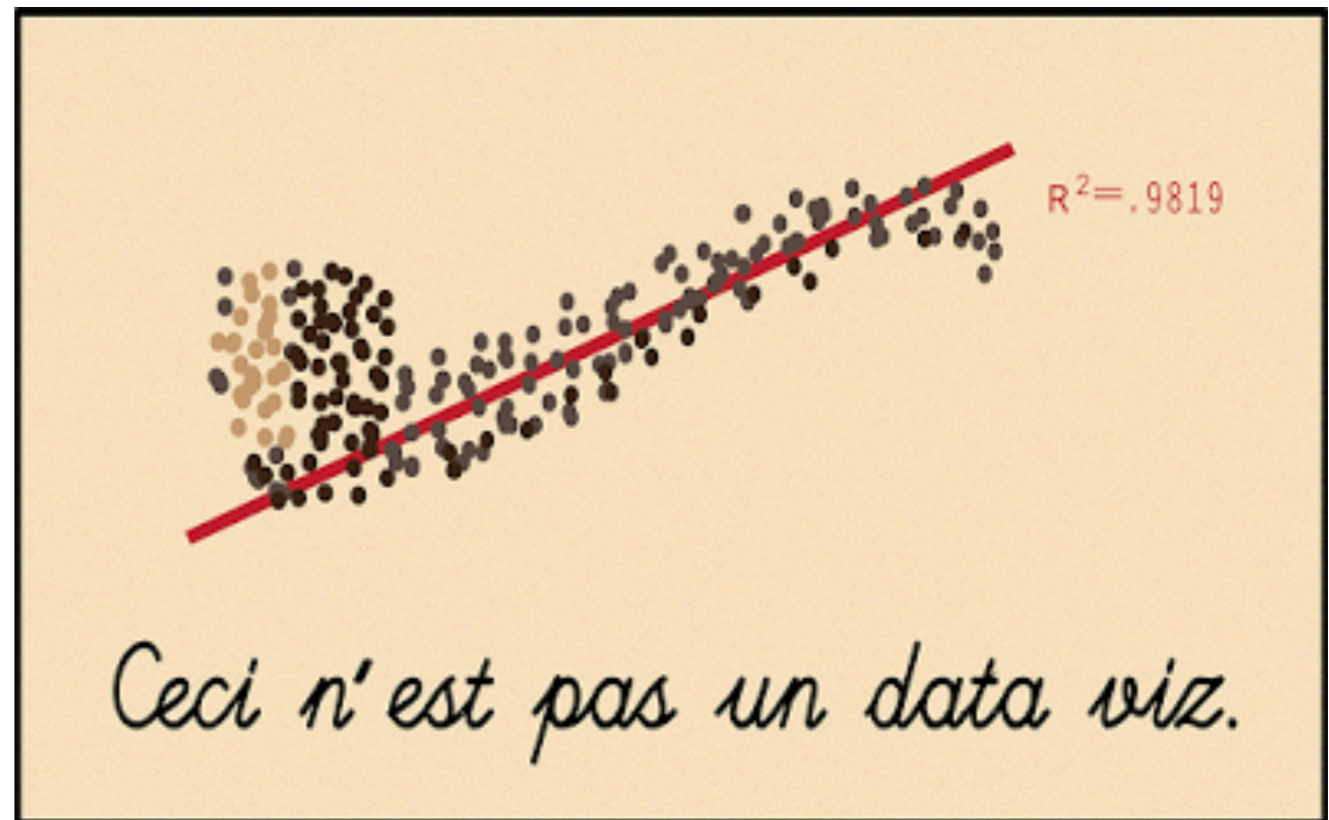
- ▶ más betöltési stratégiák
- ▶ teljesítménynövelés nagy terhelés alatt
- ▶ erőforrás csoportok a felhasználóknak
- ▶ disztribúciók és partíciók definiálása
- ▶ valódi SQL
- ▶ a konkurrencia és a sebesség közötti egyensúly megtalálása



Riportok és vizualizáció

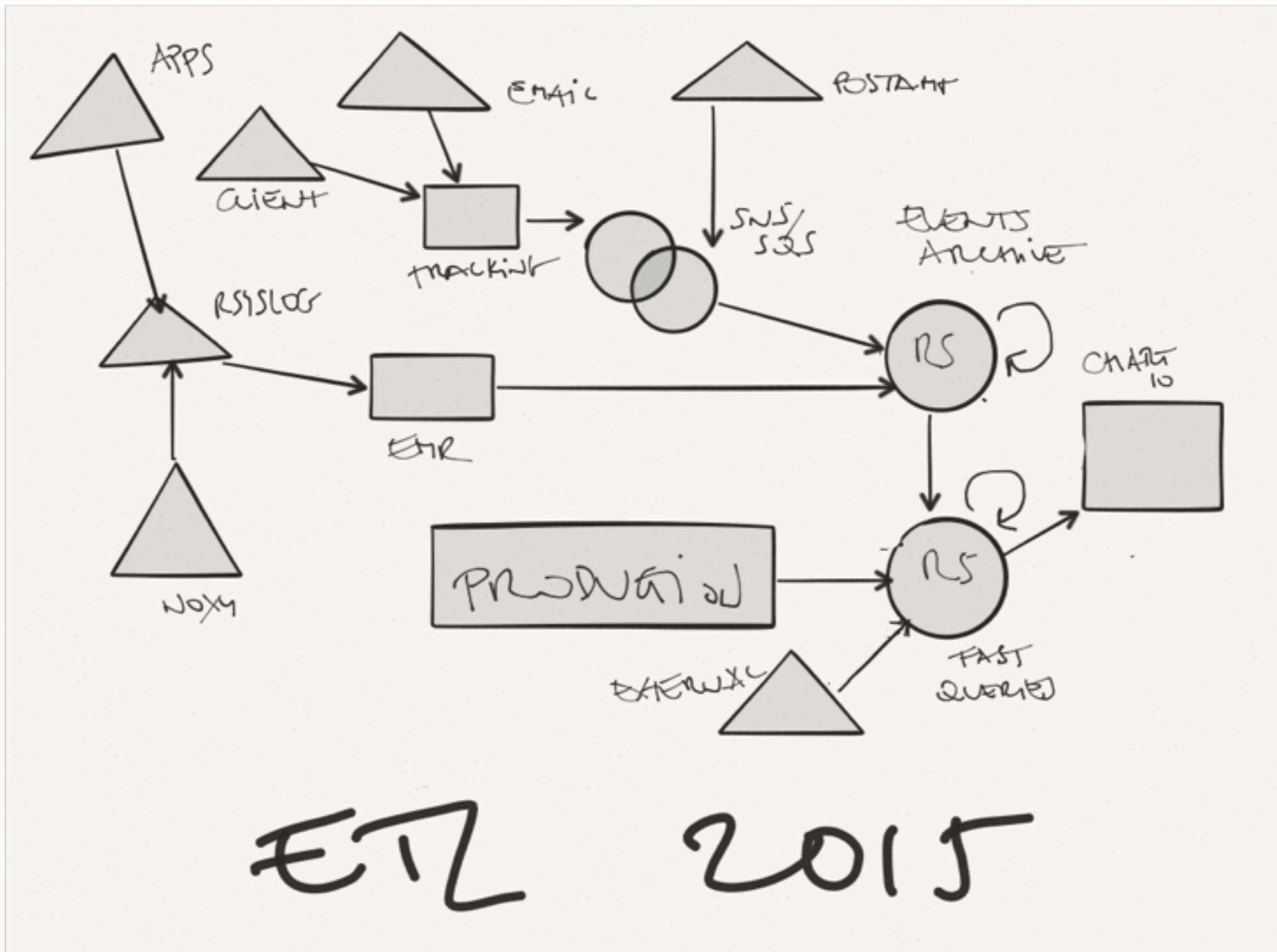
“Ha közelről nézed, igen nehéz a rossz adatot a bolondoktól megkülönböztetni.”

- ▶ Mondj nemet! Localytics, Looker
- ▶ Eszközök: Sinatra + D3, Tableau, chart.io > Power BI
- ▶ Tableau
 - + jó ár/érték arány
 - cashcow, Windows szerver, Mac app, Redshift csatoló
- ▶ 240 chart.io SQL
- ▶ Önkiszolgáló adat



Eladó, kiadó most a szívem!

Nyilas Misi pakkot kapott



Üzleti intelligencia

Jó barátok nem hagyják a másikat p-értéket számolni (ha nem értik pontosan, mi is az)

- ▶ Eszközök: Wizard (OSX), iPython notebookok
- ▶ Tipikus KPI-k, DAU (aktív), MAU
- ▶ Havi és heti kohortok
- ▶ Platform, aktivitás és földrajzi hely alapú szegmentáció
- ▶ Funnel-ek



Kísérletek

Nem vagy LinkedIn, szerencsére

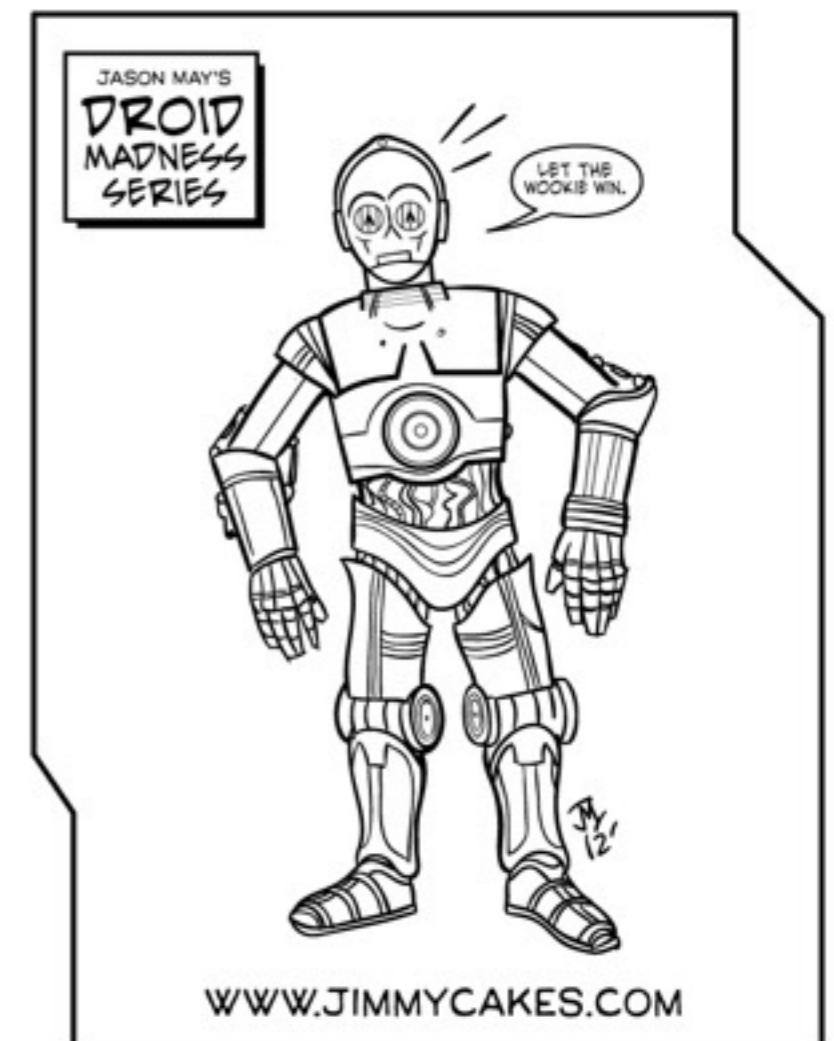
- ▶ Eszközök: Optimizely, házi fejlesztésű rendszer
- ▶ A/B tesztek funkcionalitáson és kommunikációban
- ▶ A/A – illuzórikus A/B
- ▶ túl kicsi > Bayes-i > nagyobb bizonytalanság
- ▶ rövidtávú eltérések, visszatérés az átlaghoz, véletlen variáció
- ▶ Chris Stucchio, Evan Miller



Machine Learning

PhD-vel 7 évet tölteni a Facebook-nál a javasolt ismerősök funkción

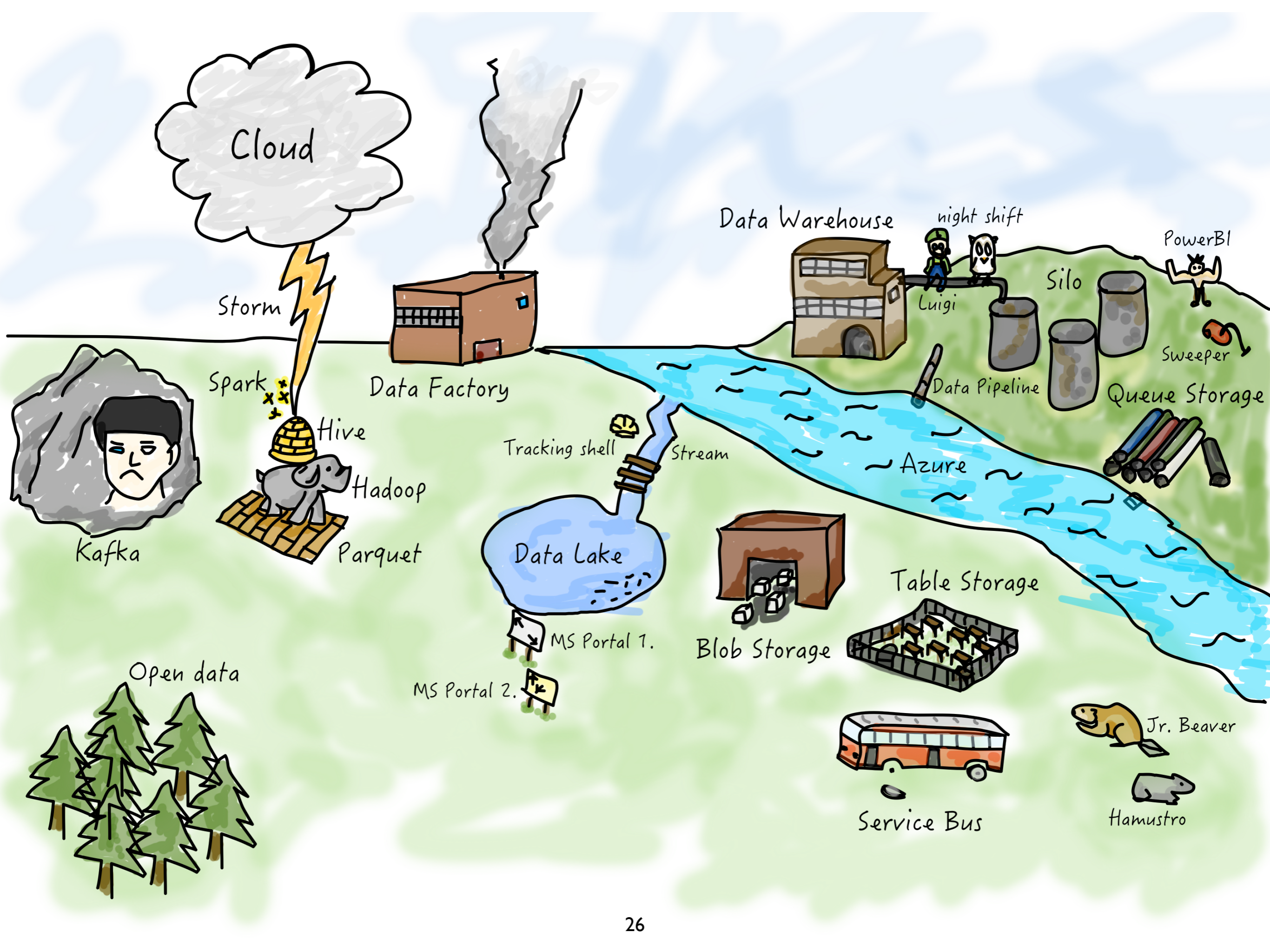
- ▶ Nincsenek PhD-seink
- ▶ Mailchimp módra
- ▶ LDA és NLP elég



Na, akkor mehet az adatinfra
AWS-ről Azure-ra?

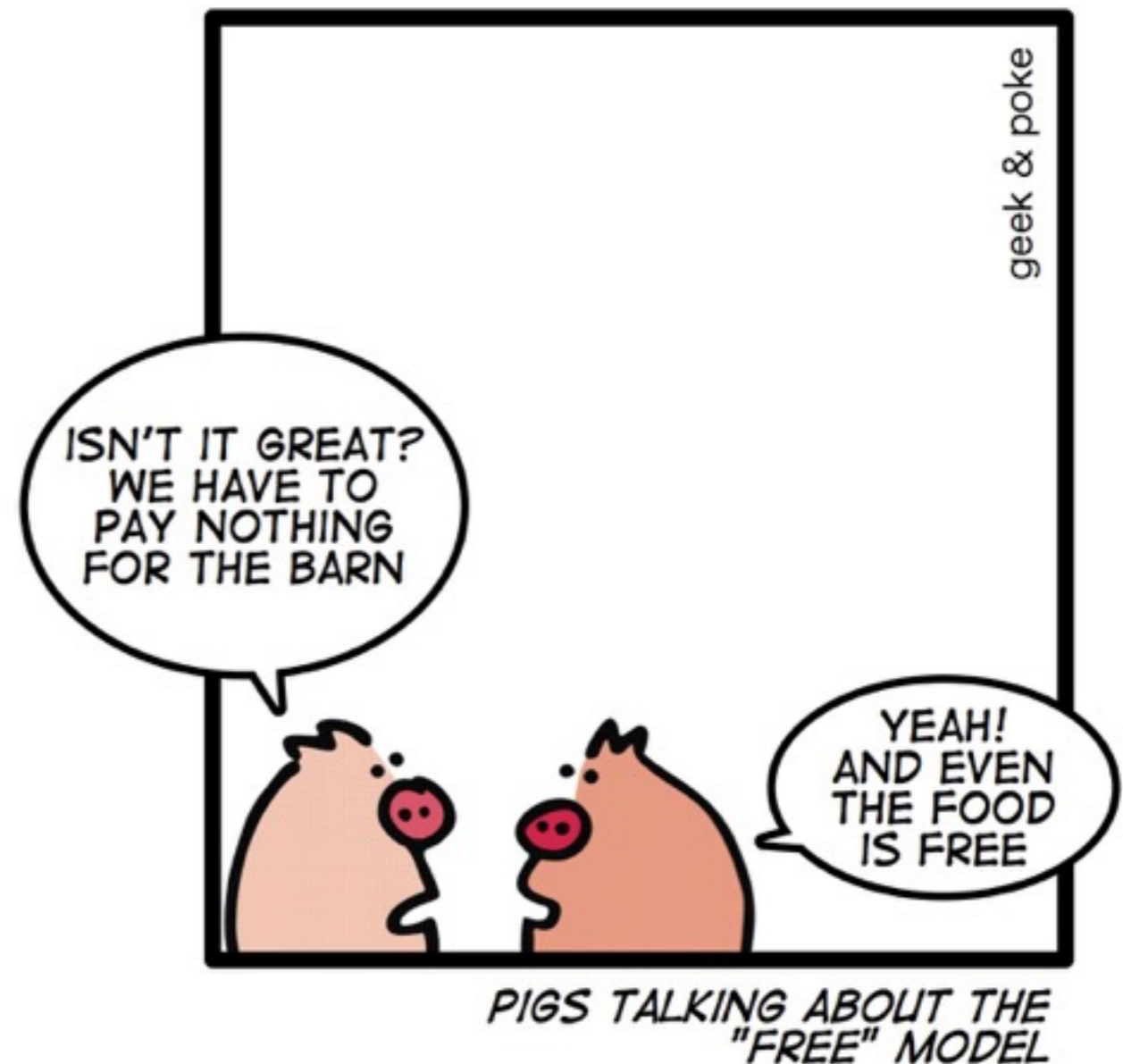
(Minden pillanatban
pontosan másfél mérnök
segítségével.)

Az adat pásztoridillje (Maciej Cegłowski)

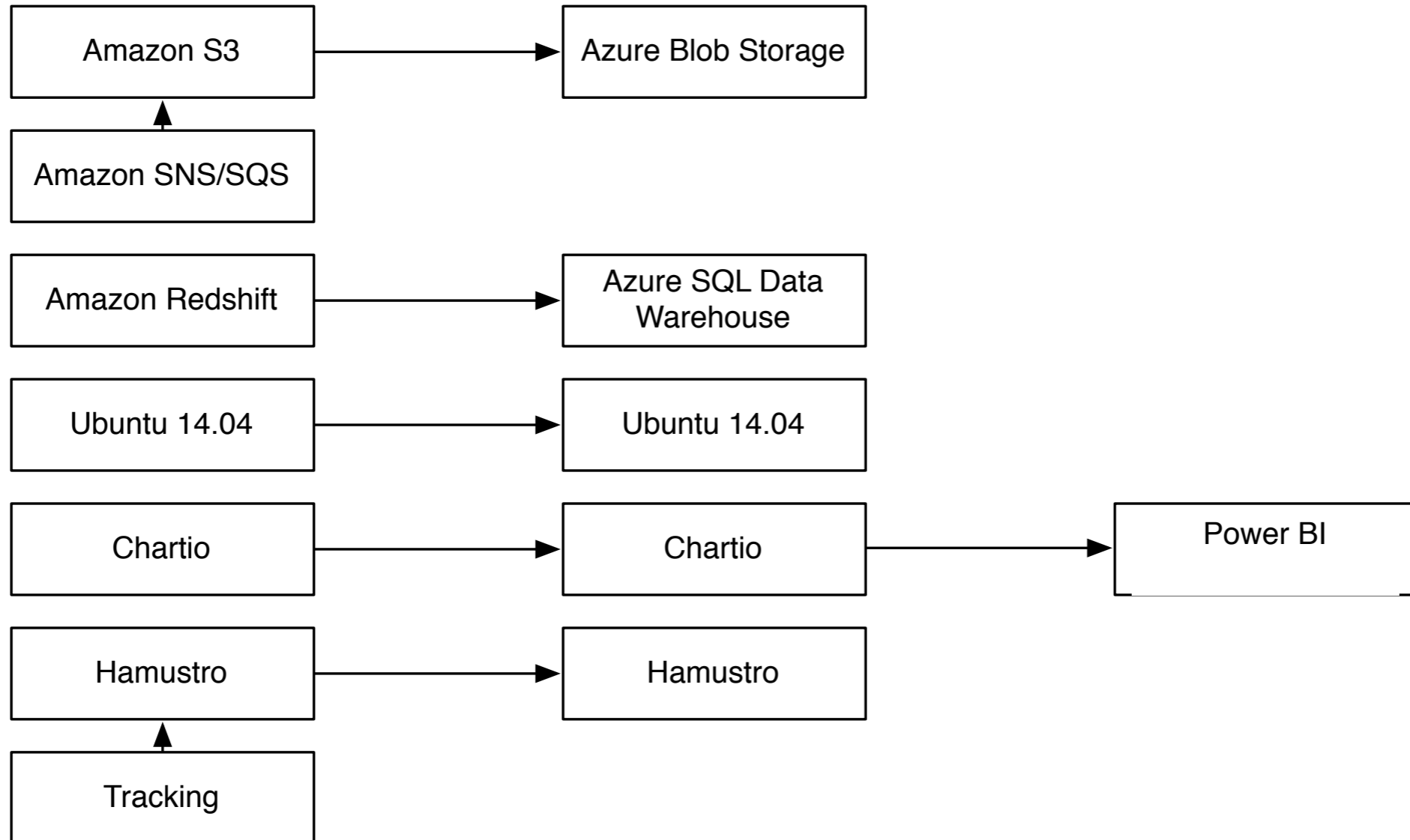


Célok

- ▶ egyszerűsítés
- ▶ az AWS specifikus részek elabsztrahálása
- ▶ felesleges komplikációk (Hadoop) eltávolítása
- ▶ Azure támogatás a komponensekben
- ▶ refaktorálás, nyílt forrásúvá tétel
- ▶ kiemelt a Total Cost of Ownership
- ▶ nem tudunk 24/7 támogatást adni
- ▶ devop monitorozás nincs fókuszbán
- ▶ ha lehet, vegyünk el belőle



Megfeleltetés és tesztmérés



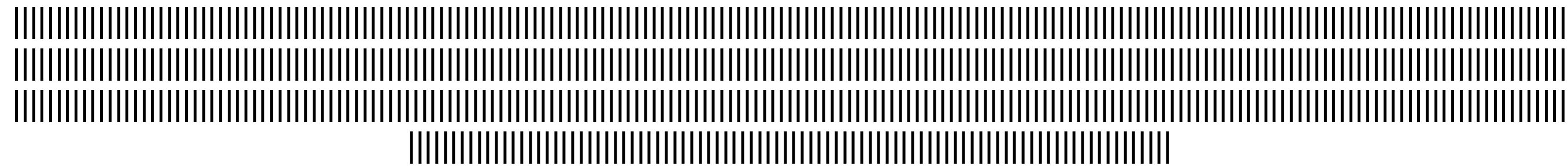
EMR-ből Jr. Beaver

- ▶ felismeri, szűri és sztenderdizálja a szerverlogokat
- ▶ térkép/csökkentő alkalmazás a felhőben
- ▶ Hadoop+Scala > make+pypy
- ▶ YAML konfigurációs fájlok
- ▶ a Pypy jobb ötlet volt, mint a Go
- ▶ nagy memória vagy nagy adat



vCPU számok

EMR (600+ db 20 számítógépben):

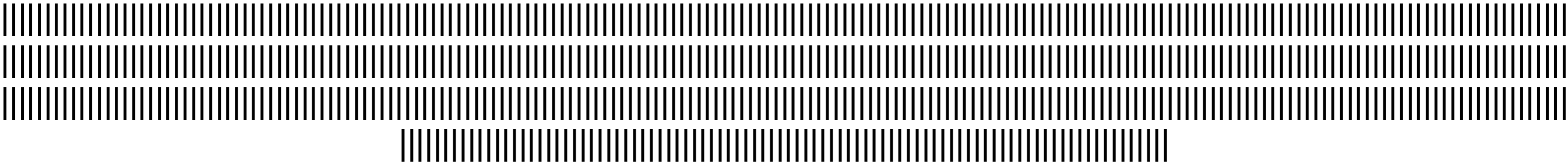


Jr. Beaver (8 db 1 számítógépben):



vCPU futásidő

EMR (600 óra):



Jr. Beaver (64 óra):



Node.js helyett

- ▶ kliensektől érkező események gyűjtése
- ▶ session és szigorú sorrendkezelés
- ▶ a felhőbe ment: Amazon SNS/SQS, Azure Queue Storage, Amazon S3, Azure Blob Storage támogatással
- ▶ percenként 6 millió esemény kezelése egyetlen 4vCPU szerverrel
- ▶ eseményformátum Protobuf vagy JSON
- ▶ Node.js helyett Go
- ▶ nyílt forrású: <https://github.com/wunderlist/hamustro>

vCPU számok

Node.js (12x1):

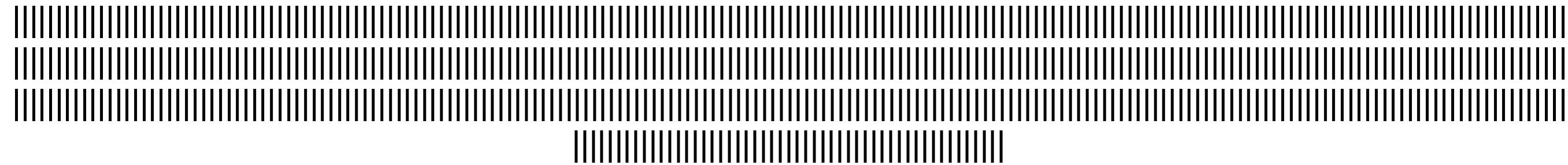
|||||

Hamustro (2x2):

|||

S3 vs. SNS 1x4CPU-val

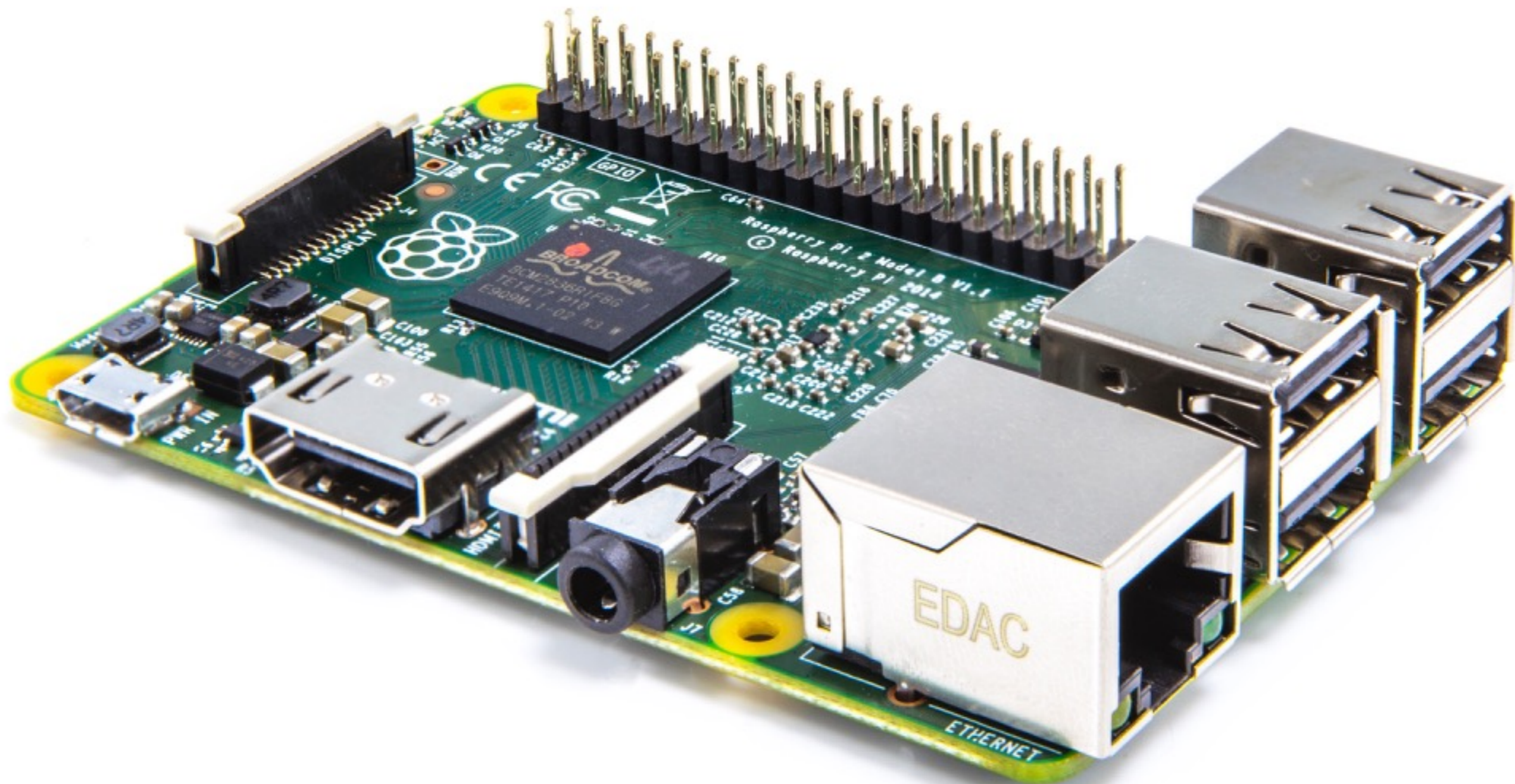
Hamustro S3 dialektusa (~6M/perc):



Hamustro SNS dialektusa (~60k/perc):



Már egy RaspberryPi3 is túlzás a jelenlegi 25k esemény/perchez



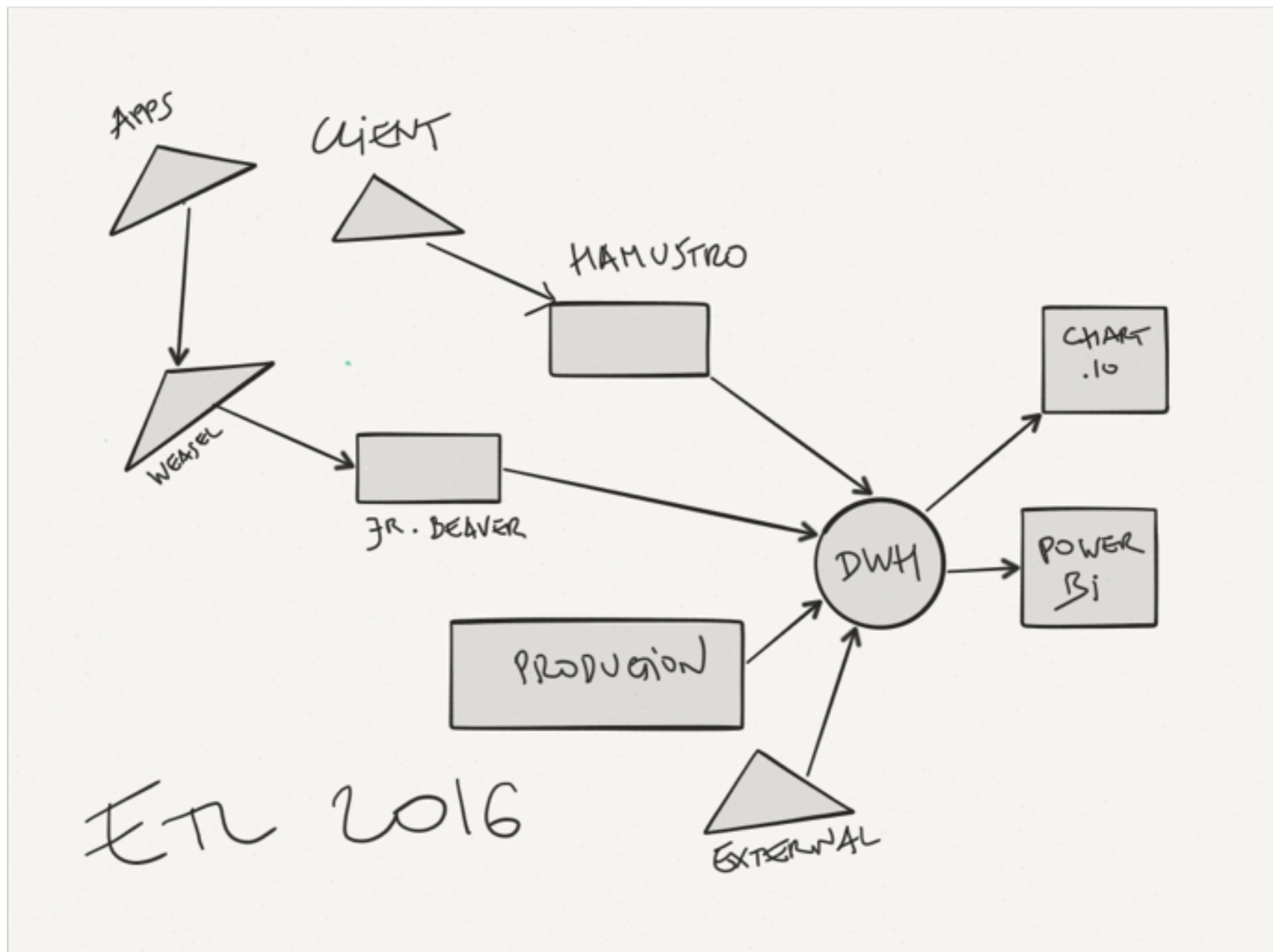
UNIX élesben Azure-on

- ▶ azrcmd <https://github.com/bfaludi/azrcmd>
CLI fájlok le- és feltöltéséhez az Azure Blob Storage-ból/be.
s3cmd-szerű funkcionalitás

- ▶ cheetah <https://github.com/wunderlist/cheetah>
CLI az MSSQL-hez OSX-re és Linux-ra.
Azure SQL Data Warehouse támogatás
psql-szerű funkcionalitás (sql-cli és sqlcmd helyett)

Tájékp csata után

Rejtélyes nyilak mindenütt



Új frázisok

hibrid, felhőagnostikus
adatinfrastruktúra

eső-után-köpönyeg (post cloud)
adatinfrastruktúra
(3 RPi összeszigszalagozva)

adatpásztor/adatkondás

Ez egy komp.

