# yql

techtalk

data processing at scale

# the plan

brief history

customers and scale

ask questions

performance dimensions

system concepts

use categories

few features

more info to explore…

# brief history

yql started as an upgrade to pipes (developers asked for cmd version of pipes)

evolved into "select * from internet" – cloud data serving

added insert update delete – grew its vocabulary
<execute> added
developer console
…lots of features

upgraded runtime architecture to scale

introduced tenancy into serving

over 1200 open tables
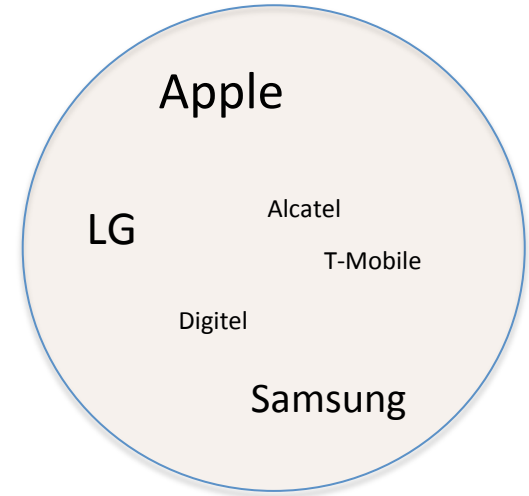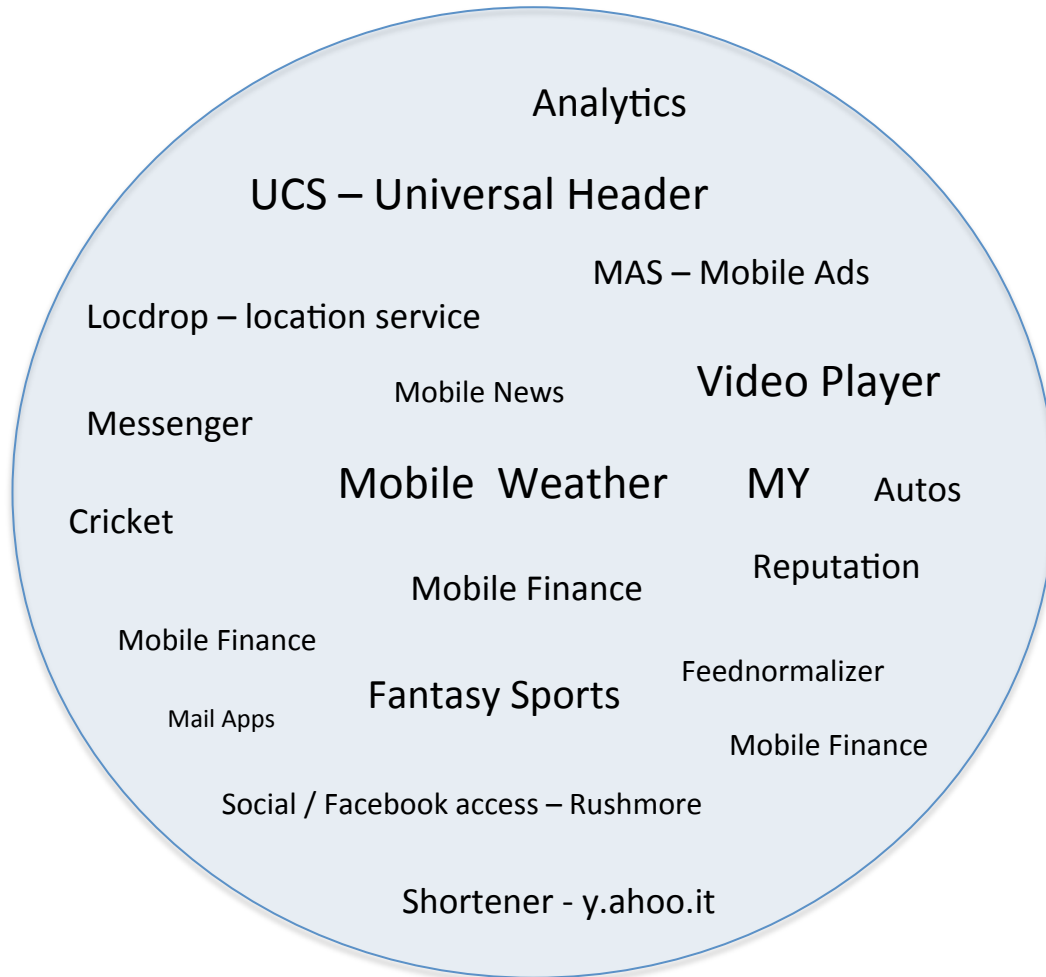over million of pipes
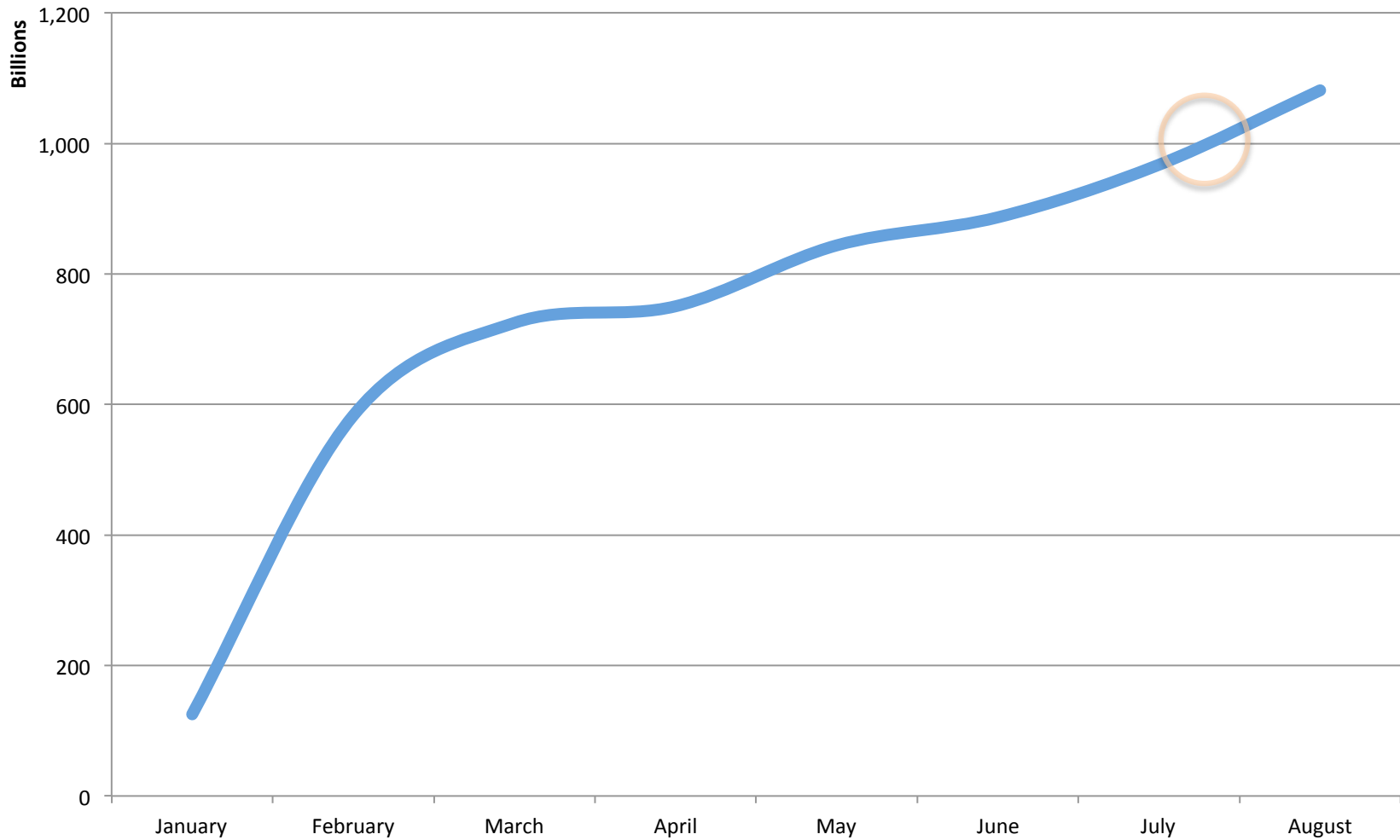hundreds of yahoo! tables added

over hundred of on-boarded customers

the mission: make data discovery, enrichment and serving fast and easy

# customers

three categories of customers

Analytics

UCS – Universal Header

MAS – Mobile Ads

Locdrop – location service

Mobile News

Video Player

Messenger

Mobile Weather    MY    Autos

Cricket

Reputation

Mobile Finance

Mobile Finance

Feednormalizer

Fantasy Sports

Mail Apps

Mobile Finance

Social / Facebook access – Rushmore

Shortener - y.ahoo.it

Apple

LG    Alcatel

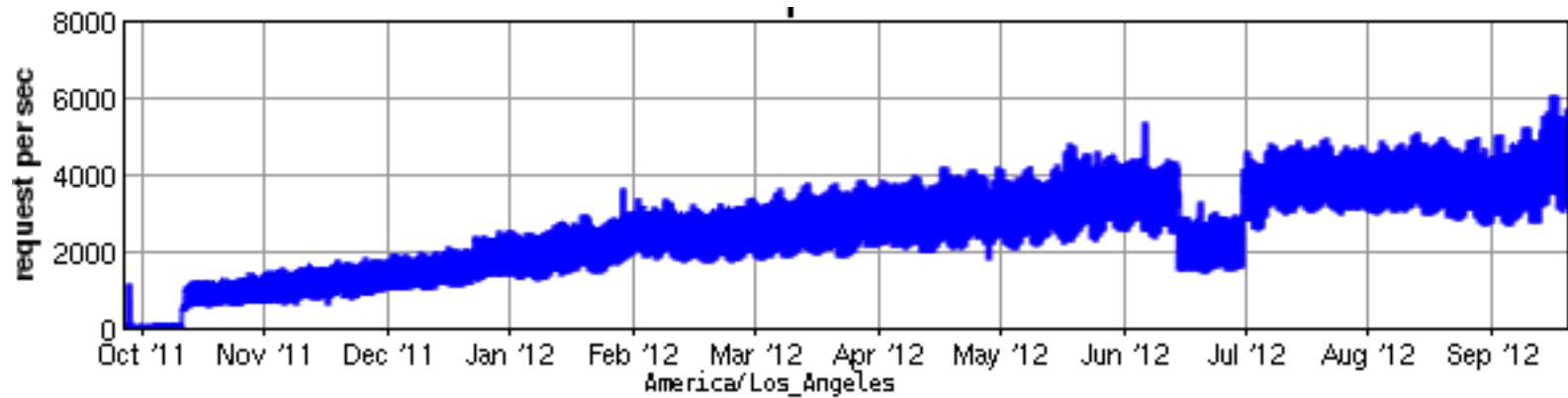T-Mobile

Digitel

Samsung

Public
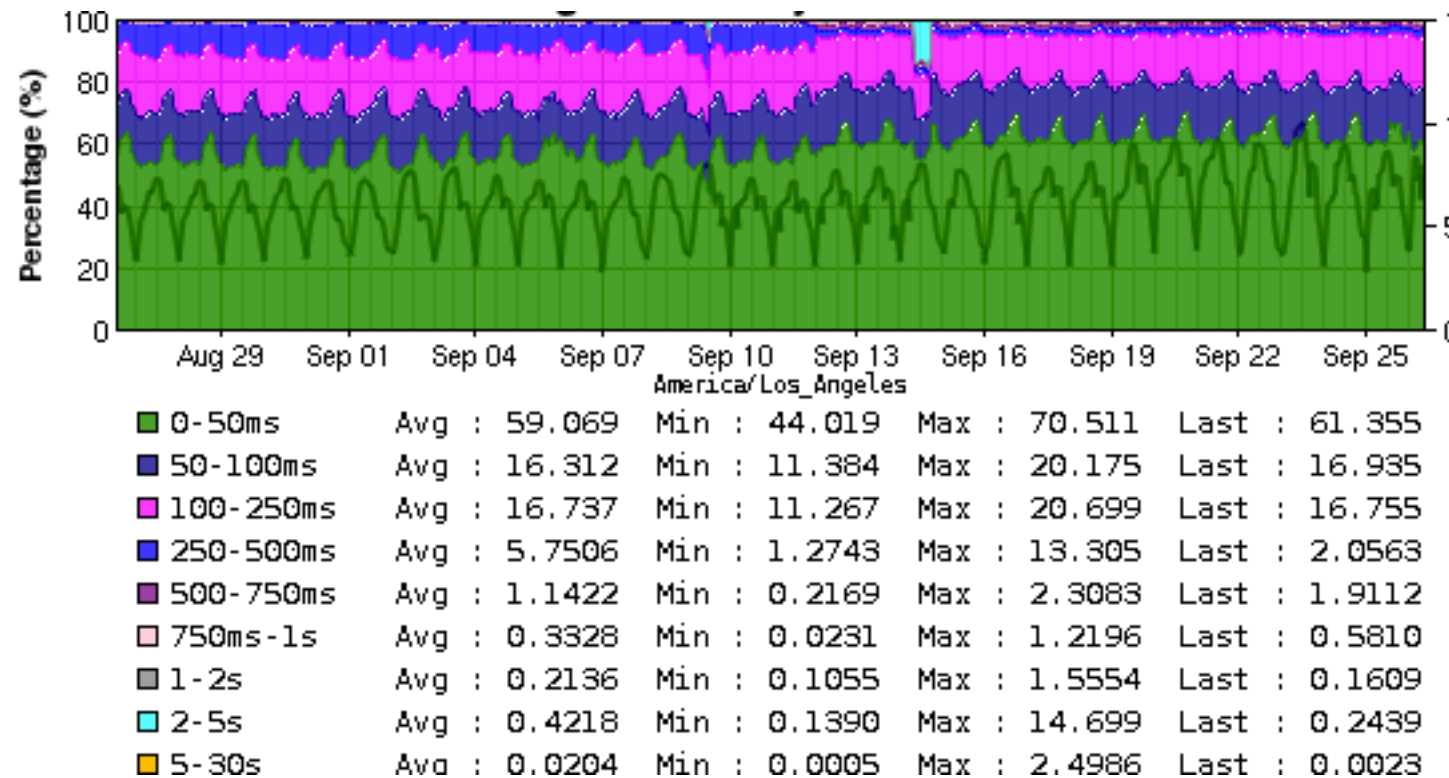
query.yahooapis.com

# yql query volume in 2012

Apple

major partner – grows with iOS adoption



predominantly weather / siri

heavily utilizes yql.multi

Apple



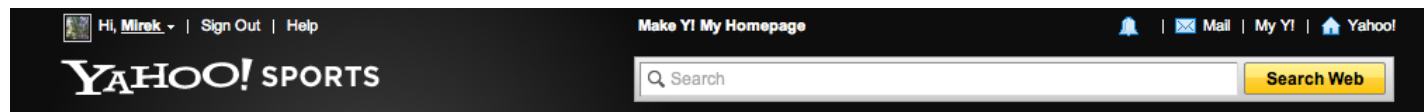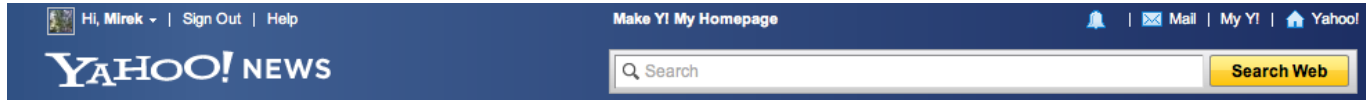| | | | | | | | | |
|---|---|---|---|---|---|---|---|---|
| 🟩 0-50ms | Avg : | 59.069 | Min : | 44.019 | Max : | 70.511 | Last : | 61.355 |
| 🟦 50-100ms | Avg : | 16.312 | Min : | 11.384 | Max : | 20.175 | Last : | 16.935 |
| 🟪 100-250ms | Avg : | 16.737 | Min : | 11.267 | Max : | 20.699 | Last : | 16.755 |
| 🟦 250-500ms | Avg : | 5.7506 | Min : | 1.2743 | Max : | 13.305 | Last : | 2.0563 |
| 🟪 500-750ms | Avg : | 1.1422 | Min : | 0.2169 | Max : | 2.3083 | Last : | 1.9112 |
| 🟫 750ms-1s | Avg : | 0.3328 | Min : | 0.0231 | Max : | 1.2196 | Last : | 0.5810 |
| ⬜ 1-2s | Avg : | 0.2136 | Min : | 0.1055 | Max : | 1.5554 | Last : | 0.1609 |
| 🟦 2-5s | Avg : | 0.4218 | Min : | 0.1390 | Max : | 14.699 | Last : | 0.2439 |
| 🟧 5-30s | Avg : | 0.0204 | Min : | 0.0005 | Max : | 2.4986 | Last : | 0.0023 |

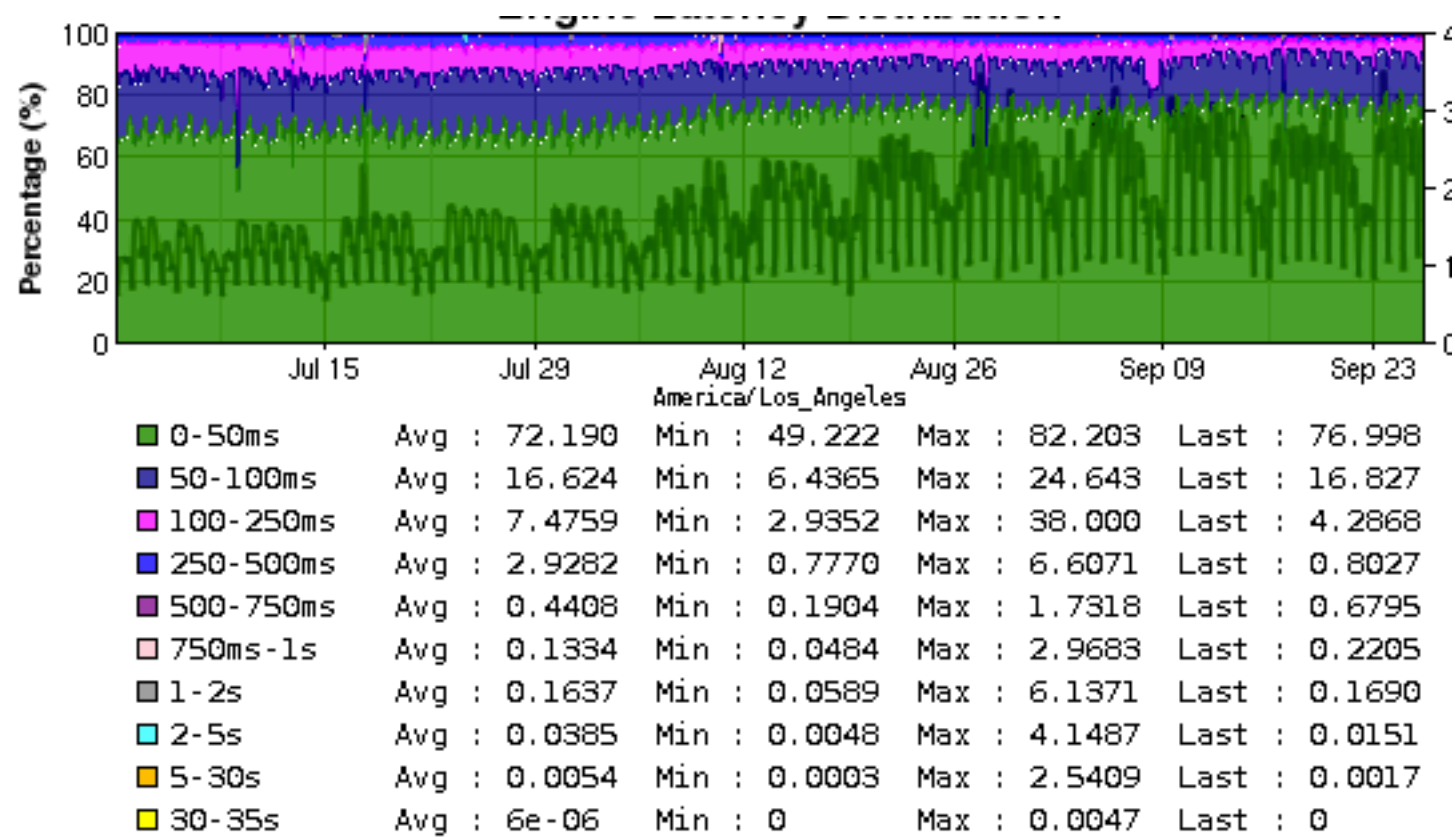places performance demands on the platform

# UCS – Universal Header



part of nearly every yahoo page

# UCS – Universal Header

# performance

lets look at platform performance as function of three dimensions

RPS

01100100

yql
QPS

01100001

**RPS**

requests per second

represents inbound demand for work and product
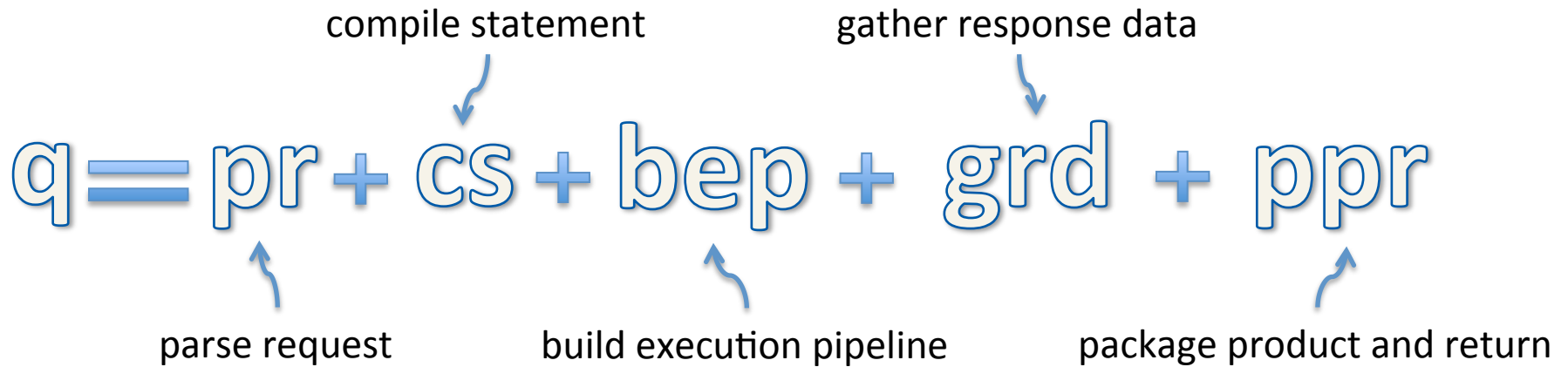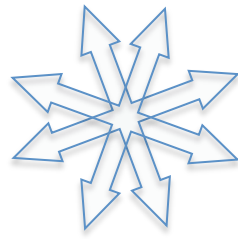
efficient system does work fast and cheap

~35K average

~60K at peak

few* platforms at y! do this much work over HTTP

what's in a yql query?
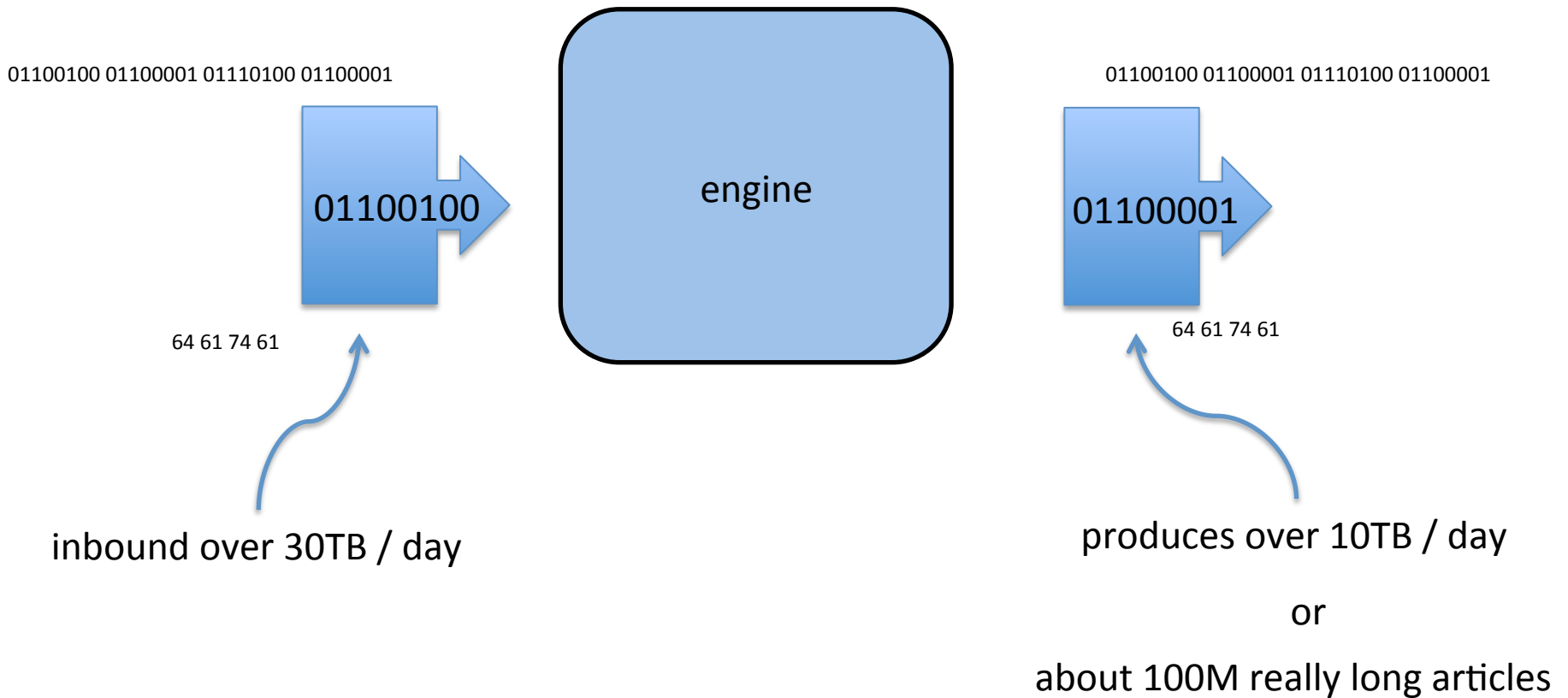
compile statement

gather response data

$$q = pr + cs + bep + grd + ppr$$

parse request

build execution pipeline

package product and return

over 400K / sec

yql queries represent work done by the engine

1T / month

01100100 01100001 01110100 01100001

01100100

64 61 74 61

engine

01100100 01100001 01110100 01100001

01100001

64 61 74 61

inbound over 30TB / day

produces over 10TB / day

or

about 100M really long articles

data volume processed by the engines

data flow by platform component per day

45TB in

in 9TB

### yts

out 14TB

out 12TB

### engine

in 31TB

out 3TB

### memcache

in 1TB

out 11TB

### http cache

in 4TB

40TB out

memcache 5%

squid 17%

yts 28%

engine 50%

data seen by the platform component

yql platform sees over 2.5 PT in a month

3B requests/day

50% < 50ms
80% < 250ms
90% < 500ms

30B+ queries/day

T1
T2
Tn

1GB / sec moved in the cloud

yql cloud = all yql system components

# system

- Deployment
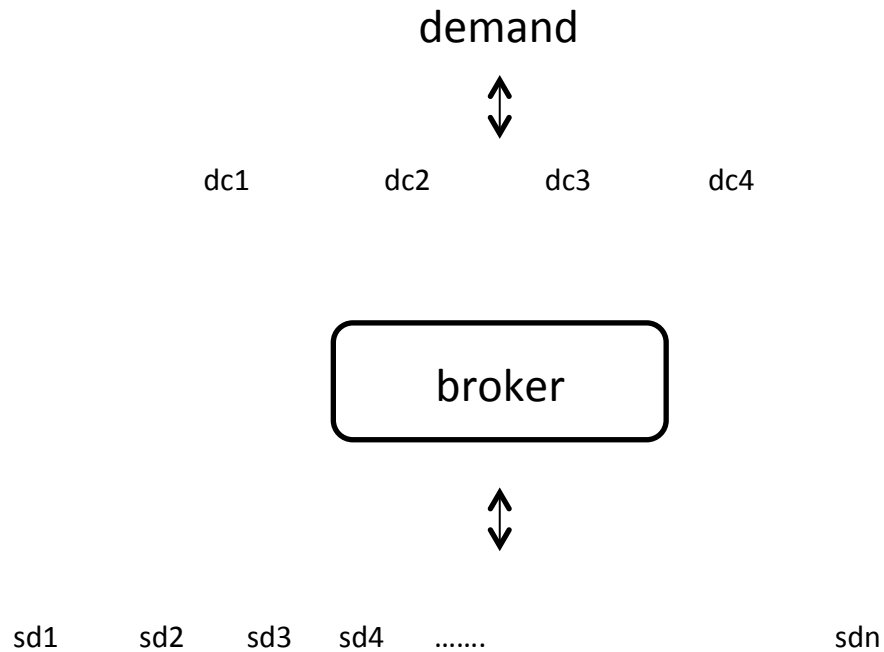  - Active in 8 colos globally
    - US: SP2, MUD, AC4, BF1 (GQ1, NE1)    = 70%
    - EU: CH1, UKL, IRD                                    = 15%
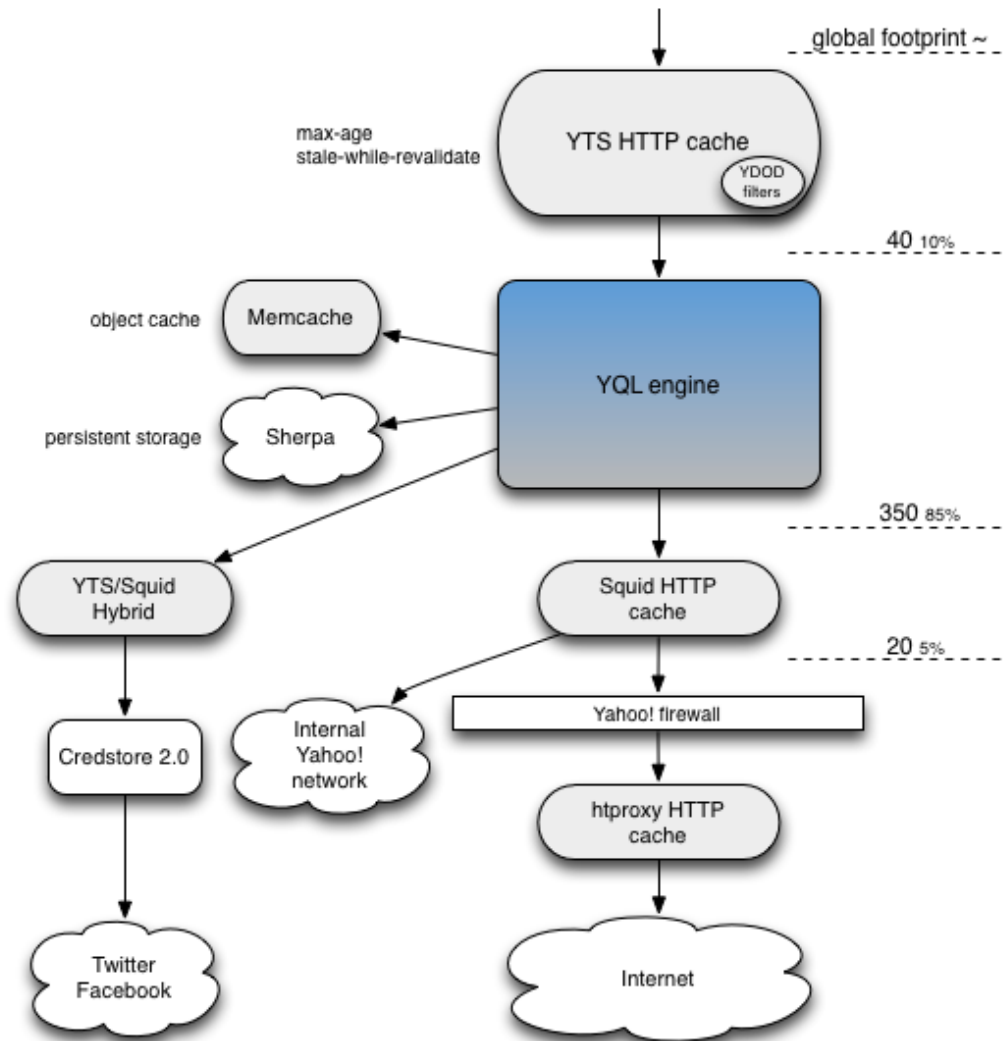    - AP: SG3                                                   = 15%

- Scales
  - horizontally
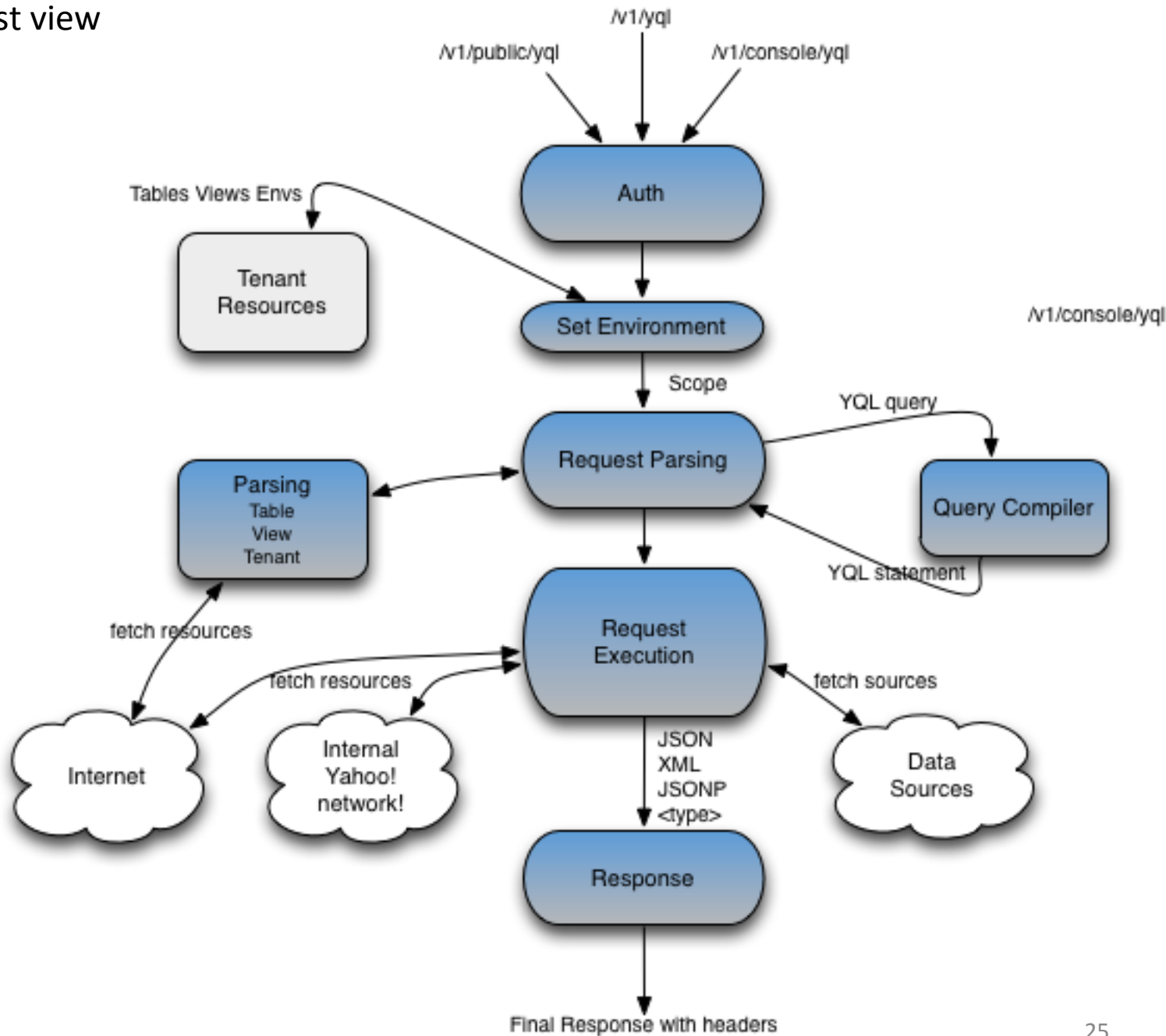  - vertically

yql platform design reflects data supply and demand model

demand

$\updownarrow$

dc1        dc2        dc3        dc4

broker

$\updownarrow$

sd1      sd2      sd3    sd4      .......                                    sdn

what does the production system look like

yql engine request view

# yql container view

uses

community

Basic                          Moderate                          Advanced

@y!

Simple – access data source and use syntax features

select * from rss where url
      in (select title from atom where url="http://y.ahoo.it/lu5F/")
            and description like "%Wall Street%" limit 10 | unique(field="title")

# Moderate – use execute to shape data

delete from twitter.status

      where id="2108869549" and username=@username and password=@password

```
<delete itemPath="" produces="XML">
    <urls>
        <url>http://twitter.com/statuses/destroy/{id}.xml</url>
    </urls>
    <inputs>
        <key id="username" type="xs:string" required="true" paramType="variable"/>
        <key id="password" type="xs:string" required="true" paramType="variable"/>
        <key id="id" type="xs:string" required="true" paramType="path"/>
    </inputs>
    <execute><![CDATA[
        y.include("http://yqlblog.net/samples/base64.js");
        var authheader = "Basic "+Base64.encode(username+":"+password);
        response.object = request.header("Authorization",authheader).del().response;
    ]]></execute>
</delete>
```

## Advanced – native implementation of data API

insert into yql.storage.admin (url) values ("http://localhost/table.xml")

not a use

caching only

(if you use yql for its cache then something went seriously wrong elsewhere in your design and implementation)

# features

select
insert
update
delete

https
streaming
charset controls
encoding

fxible creation, on demand load and reload of data application

projections

sub-select
join
y.query

internal table – extended schema
desc
xpath

yca

pipe
table
view
y.rest

oauth
developer console
proxy selection

executable business logic in Java or JavaScript
y.include

local and remote filter
authorization and access

cookies
backposting

meta services (JSON, XML, JSONP, JSONPX)
cross tenant calling

cache controls

environments

hosted
paging
flexible tenant configuration
batching

store://
easy to integrate with your development and testing flow

query, filter and combine data across Yahoo! and beyond

functions
sds

http endpoint table load
fast
easy
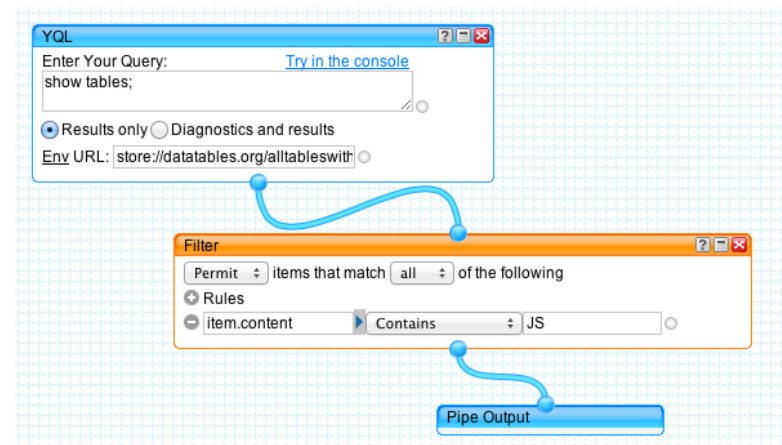async cache refresh

master those

pipe table view

pipe

a workflow
encapsulation of multiple tables and queries
expressed through JSON definition
created through GUI

instruction set bundle with operators for yql engine
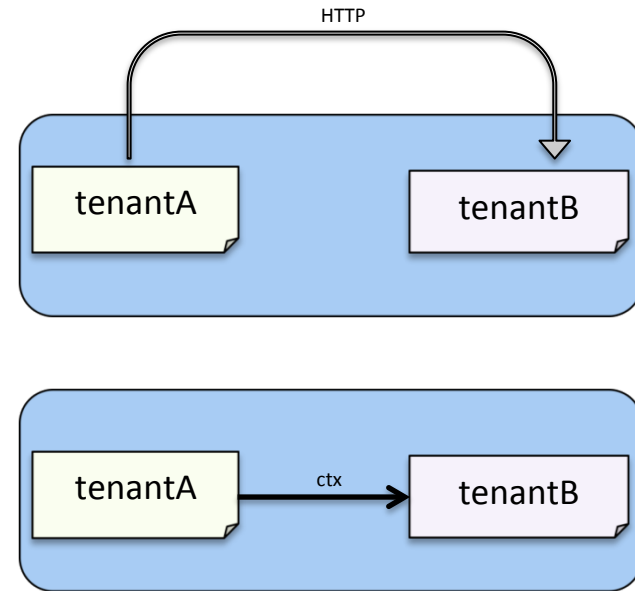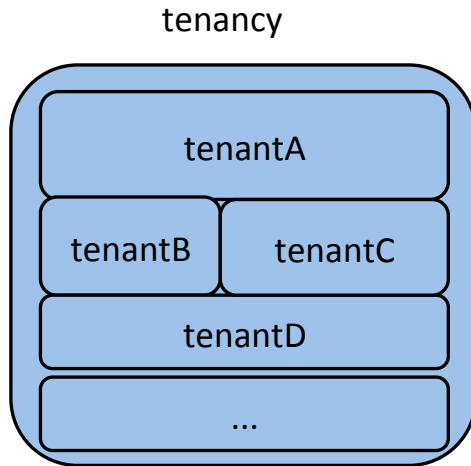
available from execute y.pipe("id")

table

yql data app definition / plugin


runtime instruction container for yql


fundamental way to express source binding

# cross tenant calling

tenancy



- keeps client resources isolated from each other
  - Tables, Views, Envs
  - Secrets (keydb)
  - Certificates (yca)
  - Customized OC TTL
- authorization schemes: public, oauth, cookie/crumb, YQL token
- custom hostnames: **media.query.yahoo.com**
- tables access can be restricted per tenant
- generate default request characteristics (b-cookie, user-agent, cache control)
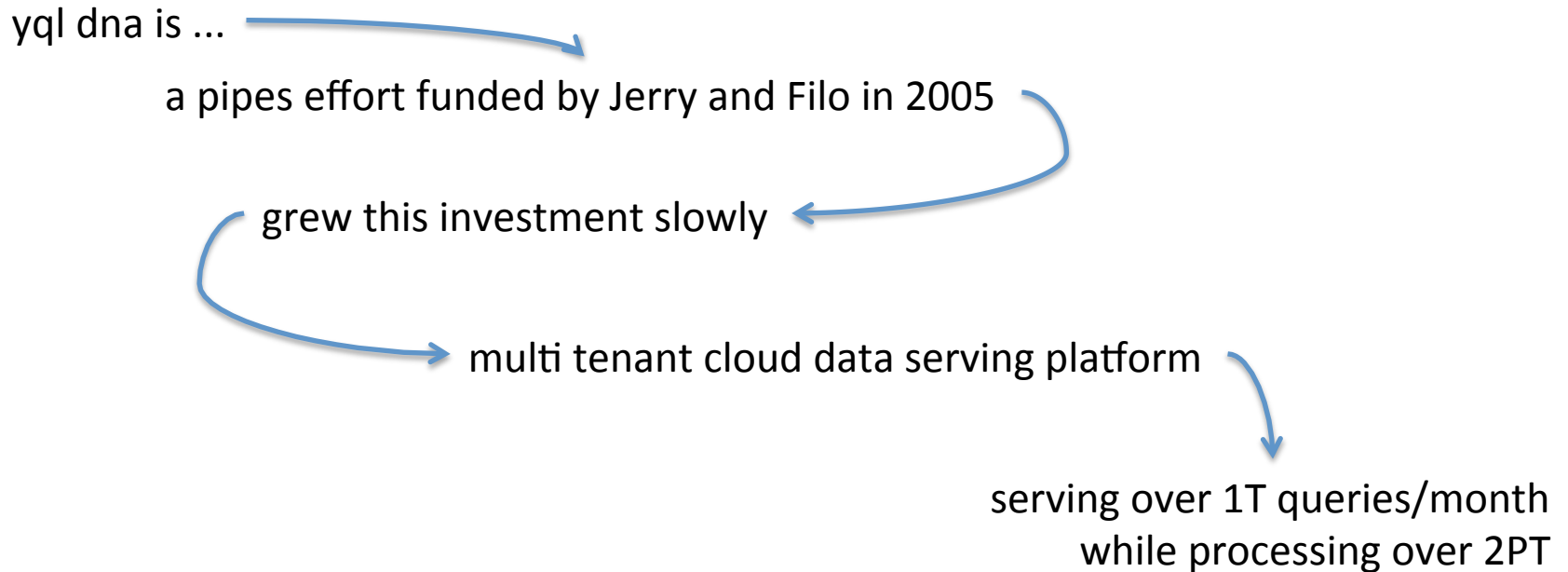
# coming soon…

views

    tenant self-provisioning

        even more flexible tenant operated resource reload

…..

    dblink

    remote table execute

# recap

yql dna is ...

a pipes effort funded by Jerry and Filo in 2005

grew this investment slowly

multi tenant cloud data serving platform

serving over 1T queries/month
while processing over 2PT

# more info...

on-boarding: <u>easy</u>

just search backyard for yql on-boarding

do <u>not</u> launch with public rate limited tenant – **query.yahooapis.com**

mirek@yahoo-inc.com

yql-discuss@yahoo-inc.com

devel.corp.yahoo.com – internal yql docs

know Java well...?  yql could be a place to write it – let us know