



emergence of micro datacenter (cloudlets/edges) for mobile computing

Victor Bahl

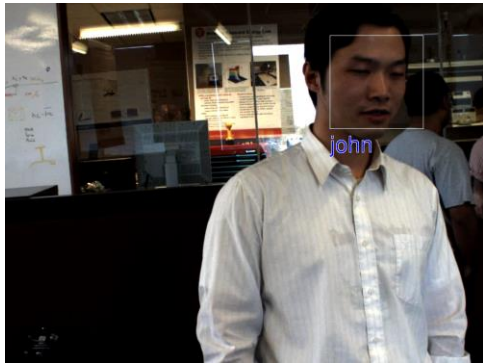
Wednesday, May 13, 2015

what if our computers could see?

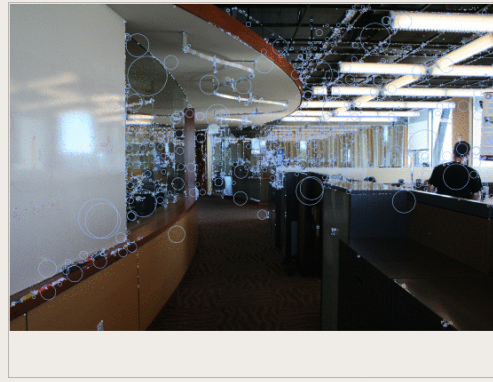


Microsoft's HoloLens

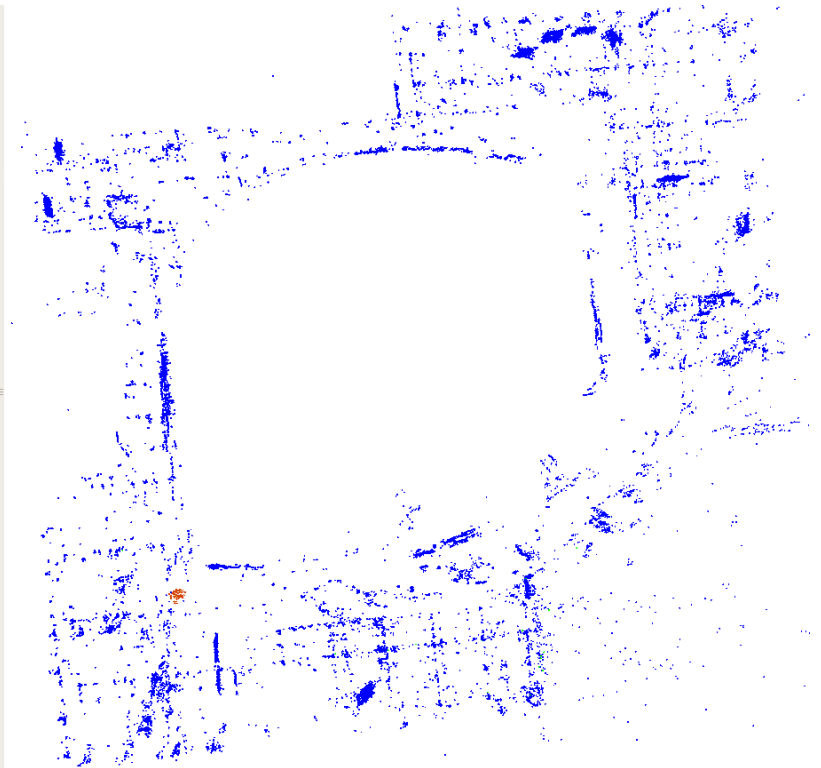
who?



where?



what?



seeing is for real

Since February of 2012, Rialto, California has required all police officers to wear a camera to monitor all interactions with the public.

After this practice was instituted, the number of complaints to the department dropped by 88%. Further, the number of instances of the police using "force" on people dropped 40%.

"When you put a camera on a police officer, they tend to behave a little better, follow the rules a little better..."
- Chief William Farrar

FB/POLICETHEPOLICEACP



BBC NEWS LONDON updated at 11:11 ET



Metropolitan Police officers start wearing body cameras

The New York Times

New York Police Officers to Start Using Body Cameras in a Pilot Program

By J. DAVID GOODMAN SEPT. 4, 2014

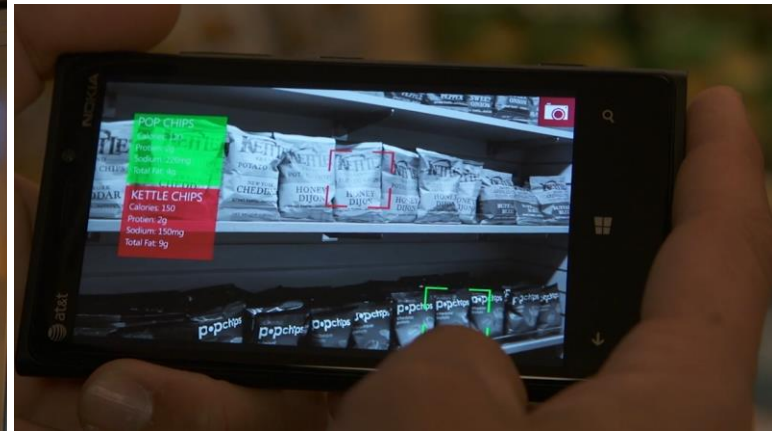
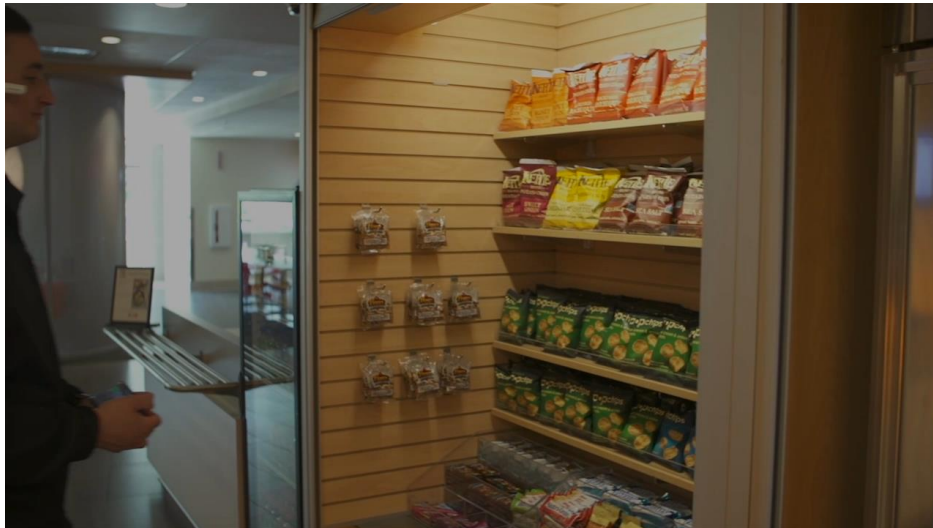
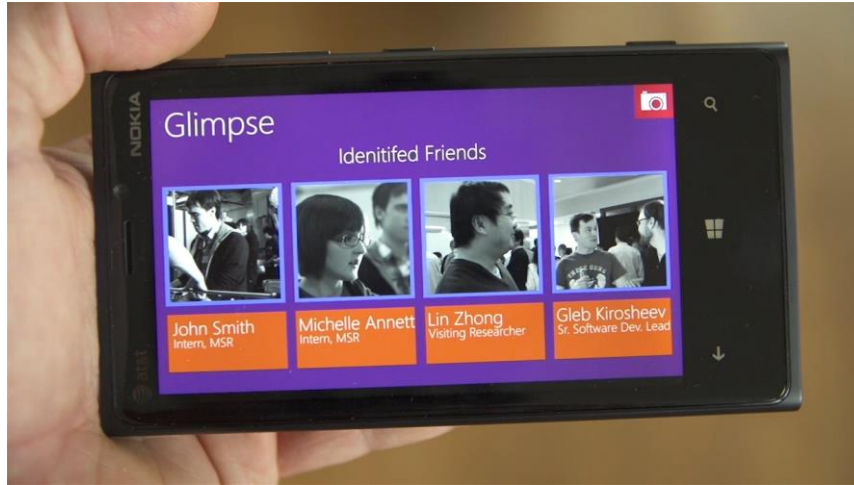
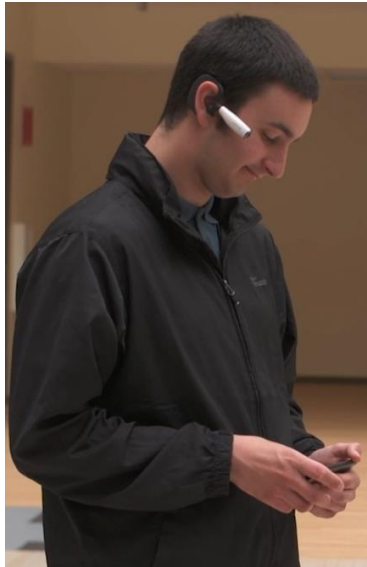


The Washington Post

D.C. police will wear body cameras as part of pilot program



MSR's Glimpse project



vision is demanding

recognition using *deep neural networks*

| | face ¹ [1] | scene [2] | object ² [3] |
|-----------------|-----------------------|-------------|-------------------------|
| memory (floats) | 103M | 76M | 138M |
| compute | 1.00 GFLOPs | 2.54 GFLOPs | 30.9 GFLOPs |
| accuracy | 97% | 51% | 94% (top 5) |

1: 4000 people; 2: 1000 objects from *ImageNet*, top 5: one of your top 5 matches

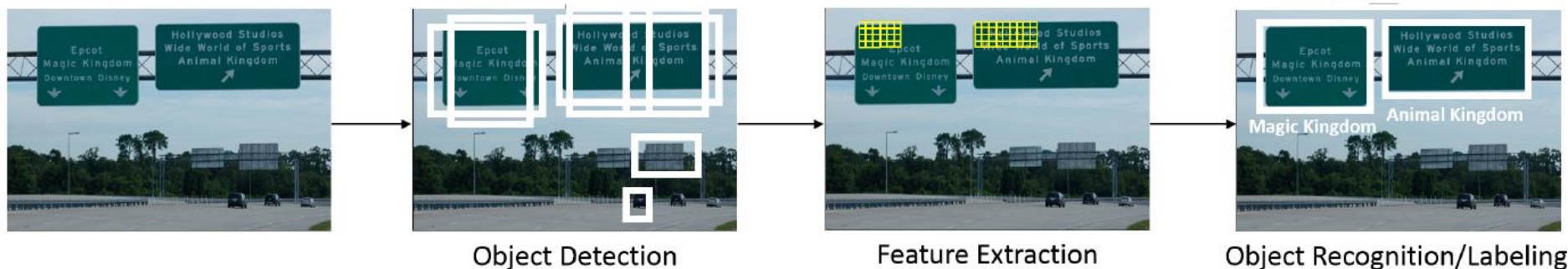
human-level accuracy, heavy resource demands
... **offloading computation is highly desirable**

[1] Y. Taigman *et al.* ***DeepFace: Closing the Gap to Human-Level Performance in Face Verification***. In CVPR 2014. (Facebook)

[2] B. Zhou *et al.* ***Learning deep features for scene recognition using places database***. In NIPS, 2014. (MIT, Princeton, ..)

[3] K. Simonyan & A. Zisserman. ***Very deep convolutional networks for large-scale image recognition***. 2014 (Google, Oxford)

recognition: server versus mobile



road sign recognition¹

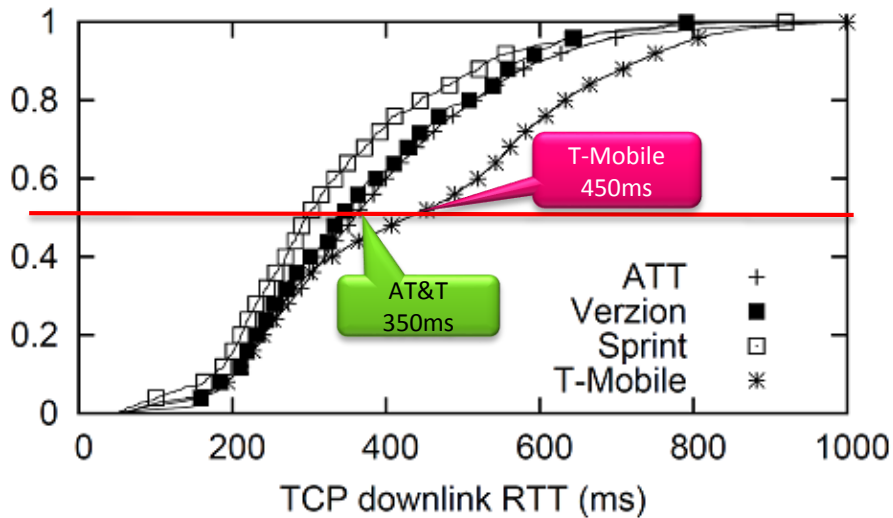
| stage | Mobile (Samsung Galaxy Nexus) | server (i7, 3.6GHz, 4-core) | Spedup (server:mobile) |
|--------------------------|----------------------------------|--------------------------------|---------------------------|
| detection | 2353 +/- 242.4 ms | 110 +/- 32.1 ms | ~15-16X |
| feature extraction | 1327.7 +/- 102.4 ms | 69 +/- 15.2 ms | ~18X |
| recognition ² | 162.1 +/- 73.2 ms | 11 +/- 1.6 ms | ~14X |
| Energy used | 11.32 Joules | 0.54 Joules | ~21X |

¹convolution neural networks

²classifying 1000 objects with 4096 features using a linear SVM

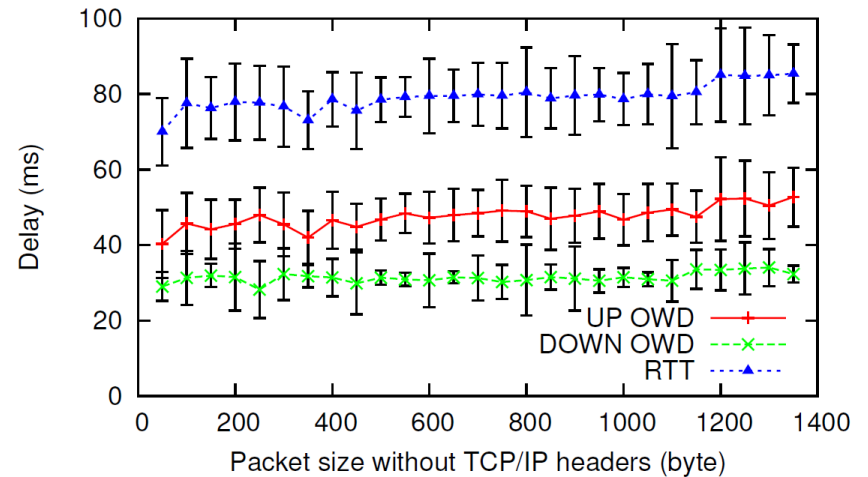
how long does it take to reach the cloud?

3g networks



MobiSys 2010

4g-lte networks



MobiSys 2013

2 years later, we still have latency issues

(May 9, 2015)

major cloud provider A

| Data Center | Average Latency |
|----------------------------|-----------------|
| West US | 115ms |
| South Central US | 131ms |
| East US | 155ms |
| North Central US | 171ms |
| North Europe | 222ms |
| West Europe | 223ms |
| Japan West | 251ms |
| East Asia | 251ms |
| Japan East | 253ms |
| Southeast Asia | 253ms |
| Central US | 276ms |
| Content Delivery Network * | 276ms |
| East US 2 | 287ms |
| Brazil South | 371ms |
| Australia Southeast | 398ms |
| Australia East | 441ms |

major cloud provider B

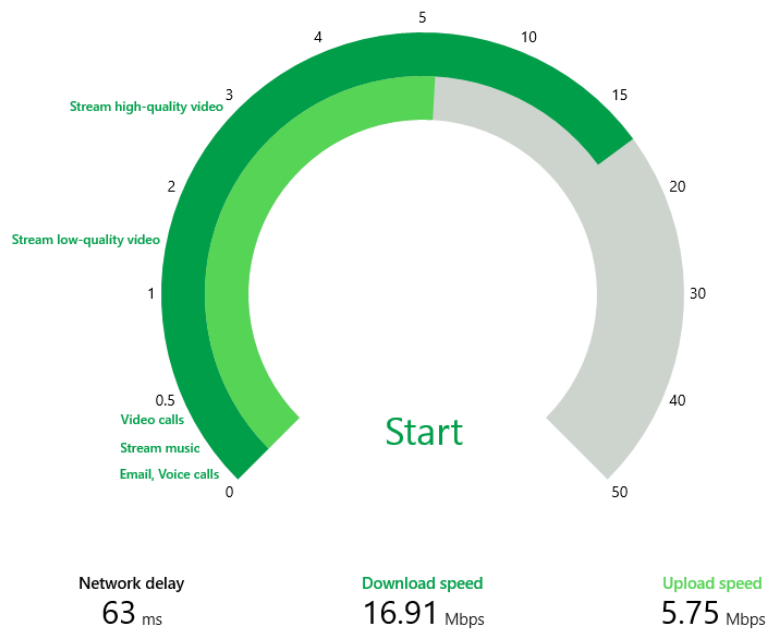
```
Command Prompt
C:\Users\bah1>tracert -d 209.85.225.99
Tracing route to 209.85.225.99 over a maximum of 30 hops
  0  0 ms  0 ms  0 ms  0.0.0.0
  1  38 ms  25 ms  47 ms  172.26.96.169
  2  45 ms  39 ms  29 ms  172.18.84.36
  3  109 ms  39 ms  39 ms  12.249.2.25
  4  59 ms  88 ms  70 ms  12.83.180.2
  5  81 ms  71 ms  88 ms  12.122.31.194
  6  76 ms  72 ms  87 ms  12.122.136.181
  7  *      *      *      Request timed out.
  8  106 ms  62 ms  80 ms  216.239.49.168
  9  81 ms  100 ms  111 ms  209.85.246.253
 10  90 ms  124 ms  112 ms  216.239.46.212
 11  110 ms  119 ms  119 ms  72.14.239.48
 12  135 ms  135 ms  133 ms  209.85.243.99
 13  138 ms  120 ms  121 ms  216.239.46.214
 14  120 ms  120 ms  125 ms  209.85.249.213
 15  192 ms  119 ms  123 ms  209.85.247.6
 16  116 ms  126 ms  112 ms  64.233.175.45
 17  *      *      *      Request timed out.
C:\Users\bah1>
```

also, <http://claudit.feld.cvut.cz/claudit/rtdata.php>

try it out – download Microsoft's Network Speed Test

Network Speed Test

Last Test (2/12/2013 12:50 PM)



Current network

Connection type
Wi-Fi
Network name
A-MSFTWLAN
Internet status
Internet access
Host name
minint-6d
Access point BSSID
6C:F3:7F:4F:88:72
Authentication
Rsn
Encryption
Ccmp



Settings

Network Speed Test
By Microsoft Research

Options

About

Privacy Statement

Contact us

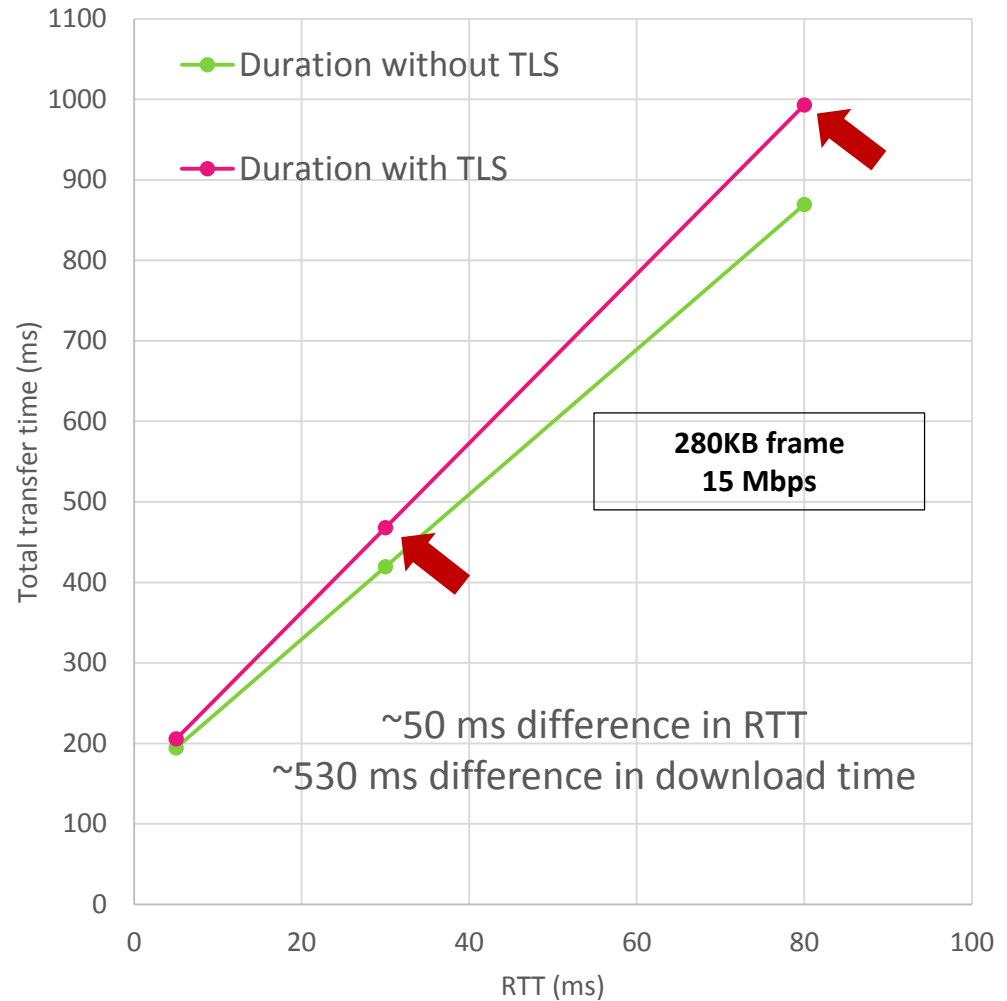
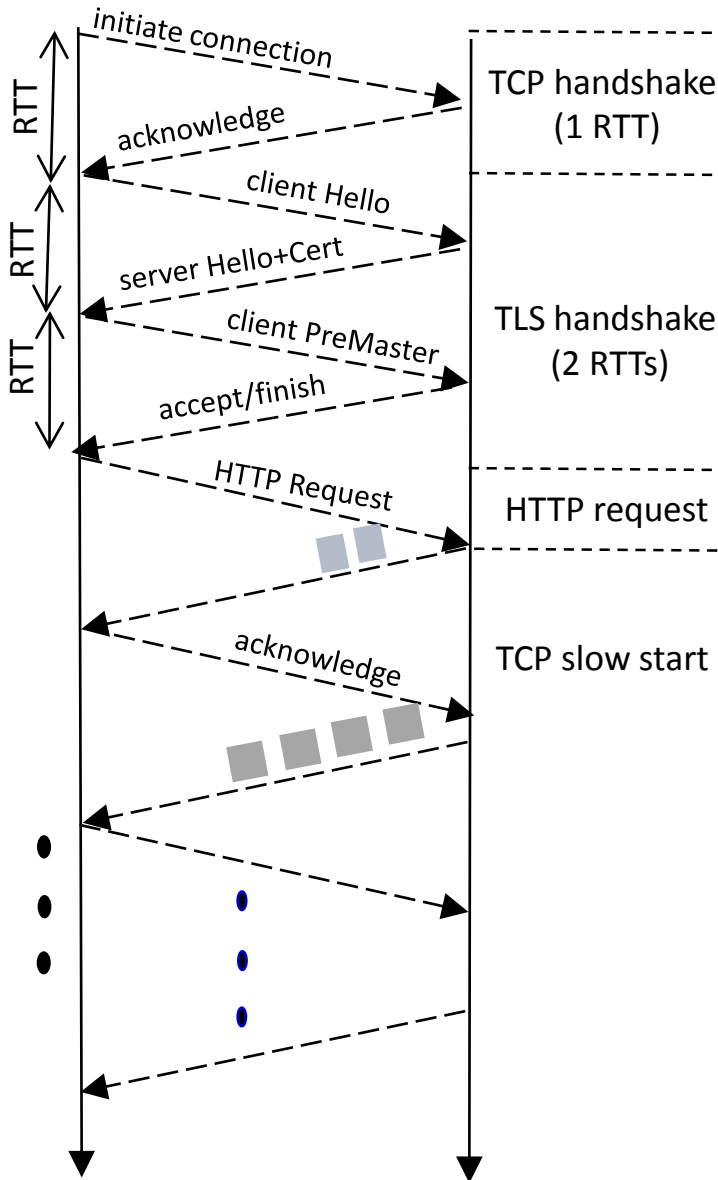
Permissions

A-MSFTWLAN
 98
 Screen
 Notifications
 Power
 Keyboard

Change PC settings

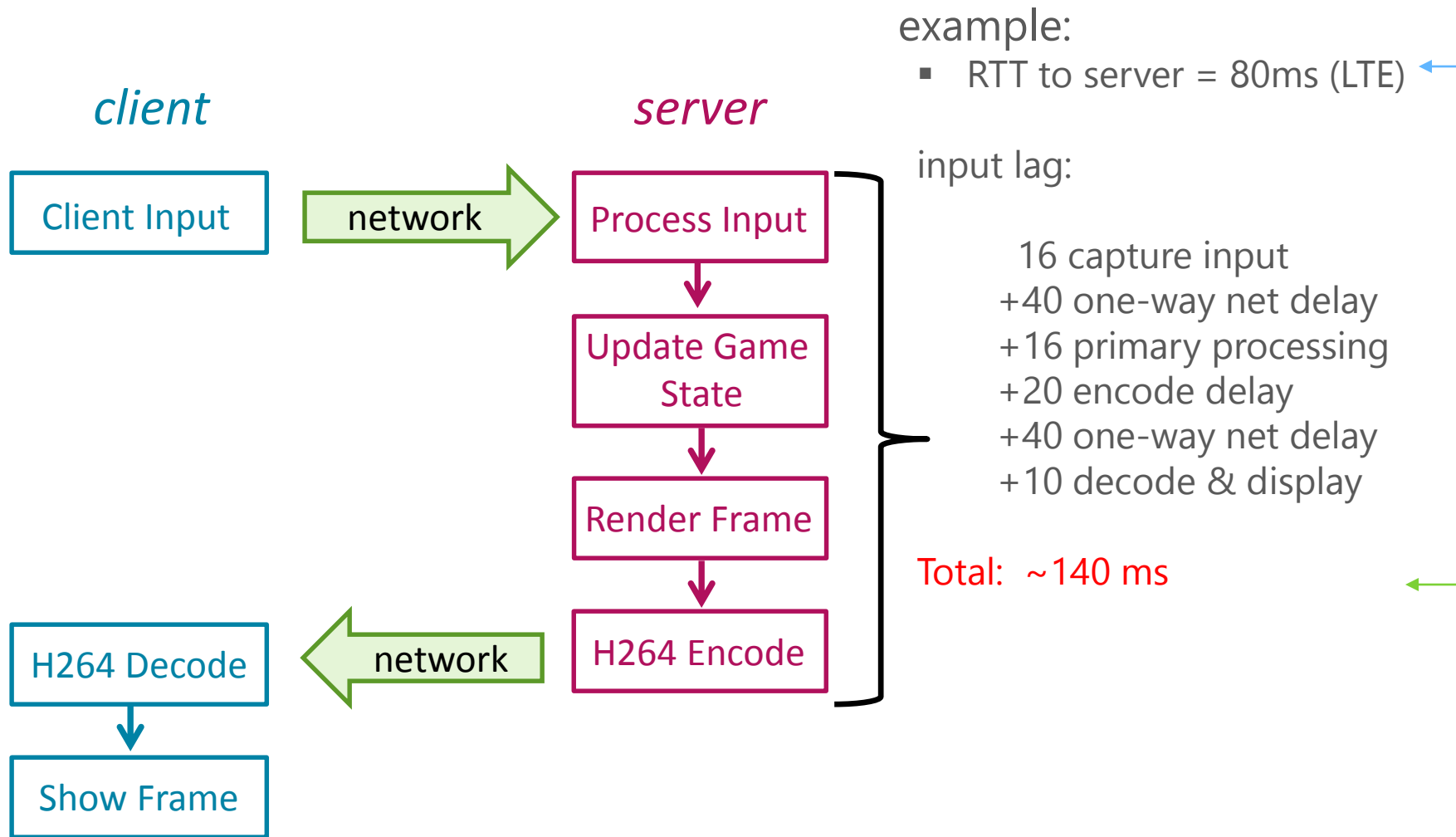
Available on Windows Phone and Windows 8

popular protocols make things worse!

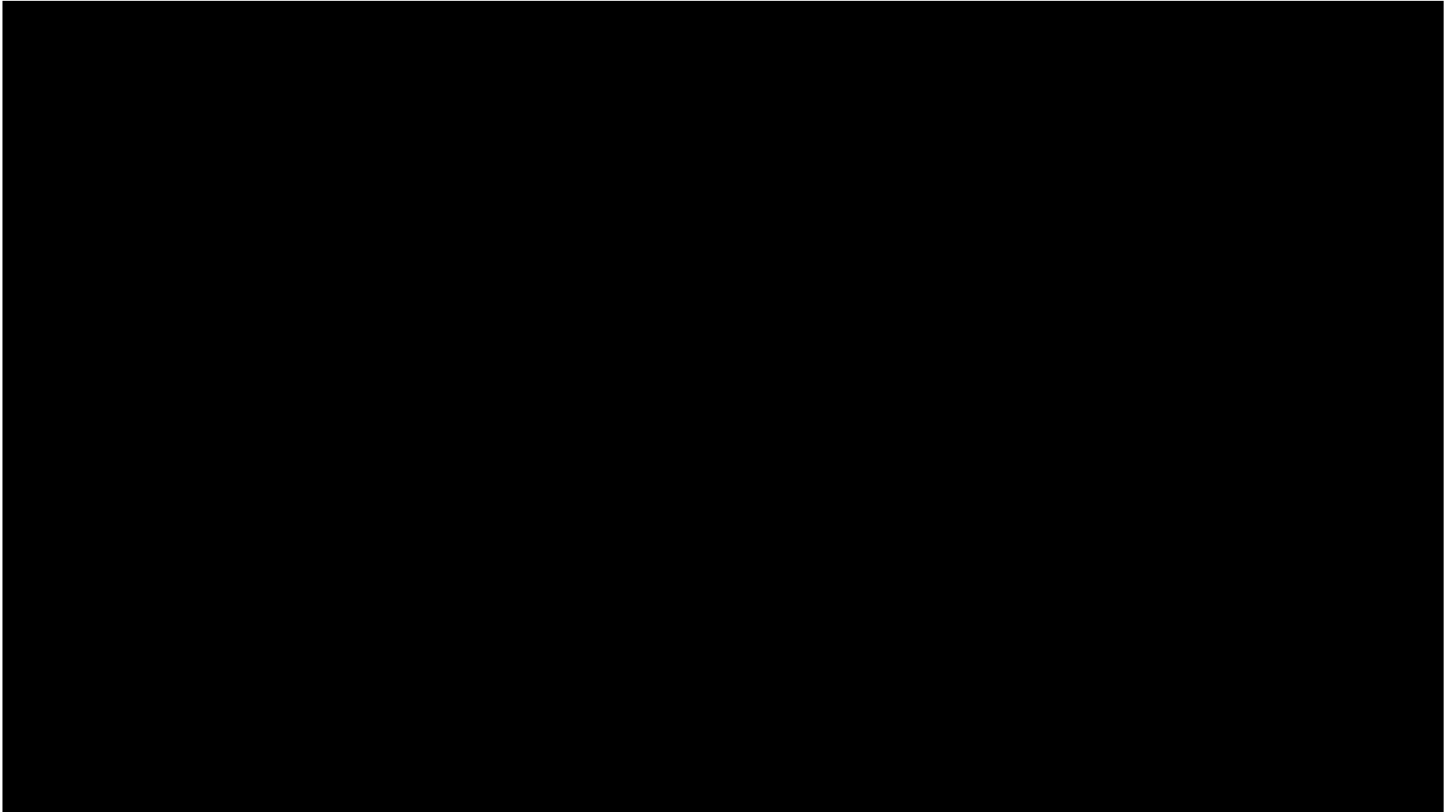


even with UDP - end user impact

fast action cloud gaming



impact of 5, 30 & 80 msec latency (fast action gaming)



latency matters!

"being fast really matters...half a second delay caused a 20% drop in traffic. and it killed user satisfaction"

- Marissa Mayer @ Web 2.0 (2008)



"...a 400 millisecond delay resulted in a -0.59% change in searches/user",
[i.e. Google would lose 8 million searches per day - they'd serve up many millions fewer online adverts]

- Jake Brutlag, Google Search (2009)



"...for Amazon every 100 ms increase in load times decreased sales with 1%"

- Andy King, book author



"...when 50% of traffic was redirected to our edges preliminary results showed a 5.9% increase in click-thru rates"

- Andy Lientz, Partner GPM, BingEdge (2013)



the fact of the matter is ...

offloading computation to a resource-rich cloud brings the true power of CS into your hands

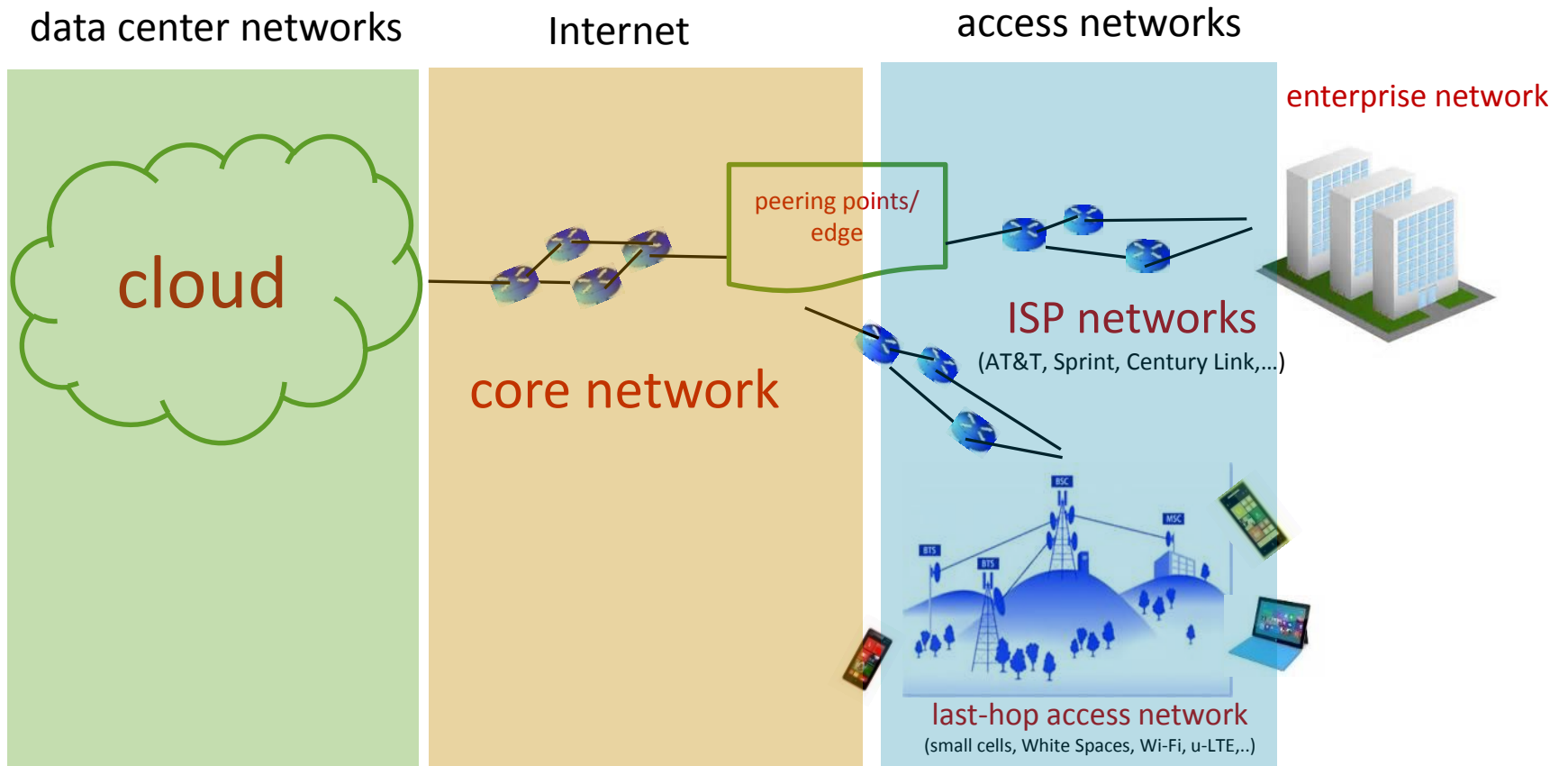
high latency & jitter to the cloud can make cloud services unusable

poor performance impacts revenue and turns users away

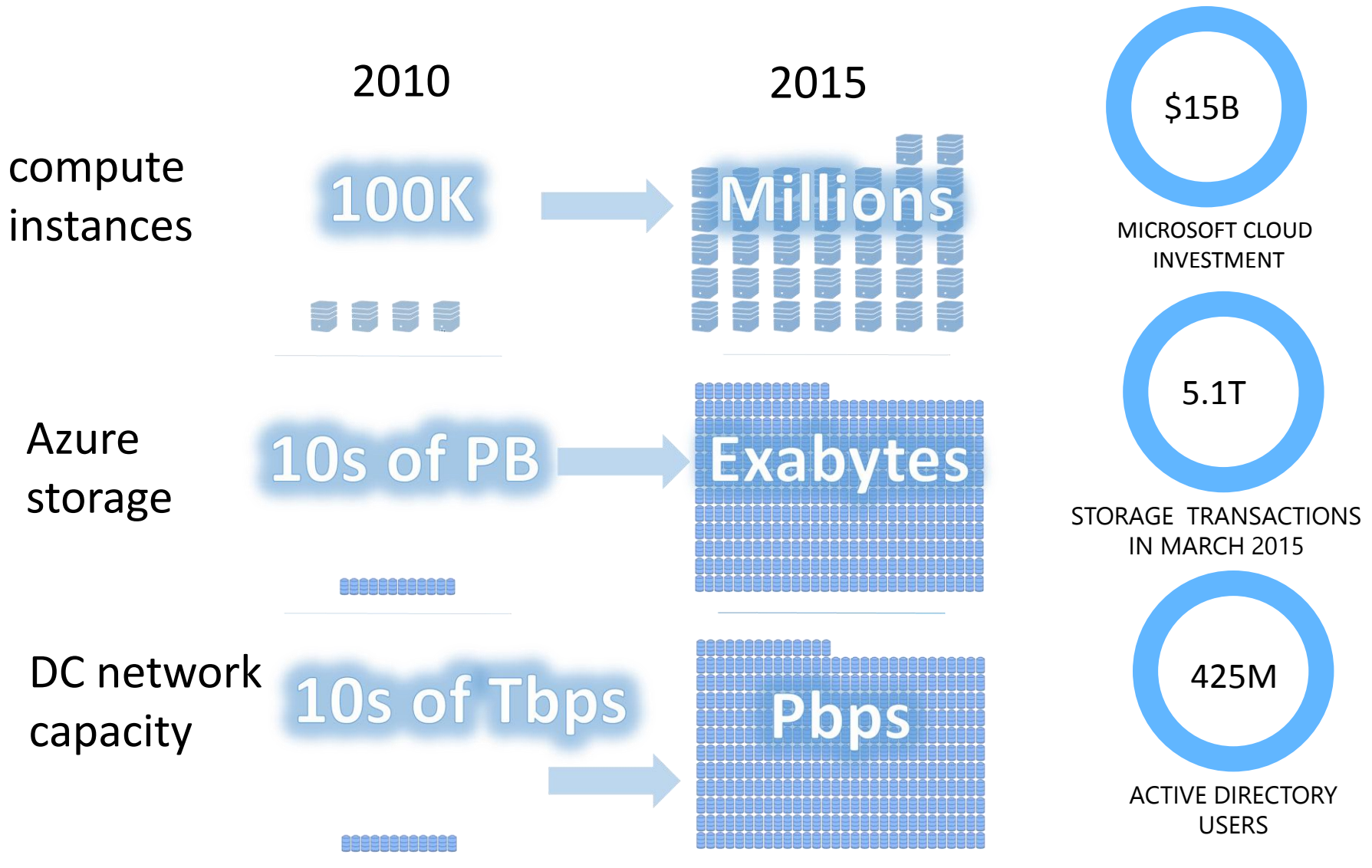
... and we have a latency problem

reducing latency

contributors to latency

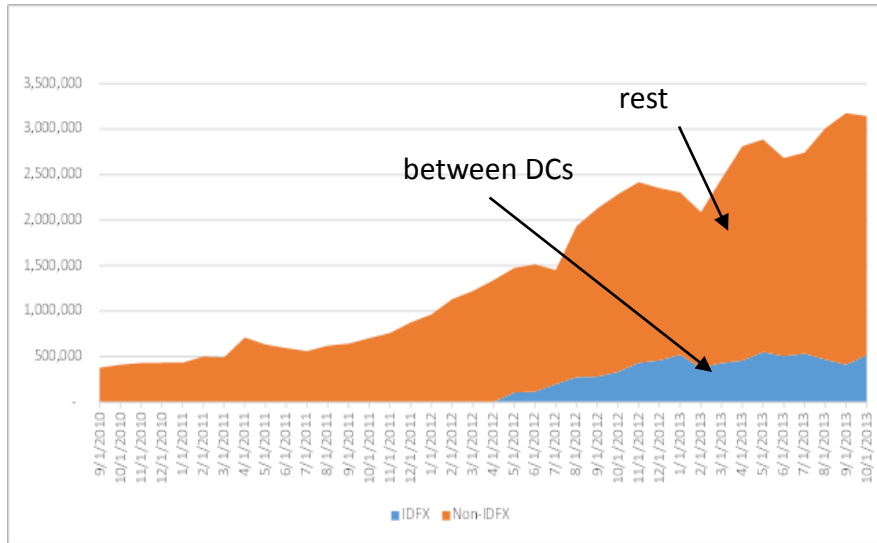


Microsoft's hyper-scale cloud



Microsoft's hyper-scale network

Microsoft's network is one of the largest in the world



1.4M

MILES OF FIBER
(DC & WAN)

4X

WRAP THE EARTH IN
NORTH AMERICAN FIBER

massive traffic growth is stressing the underlying core networks

areas MSR researchers are working on:

SIGCOM 2014

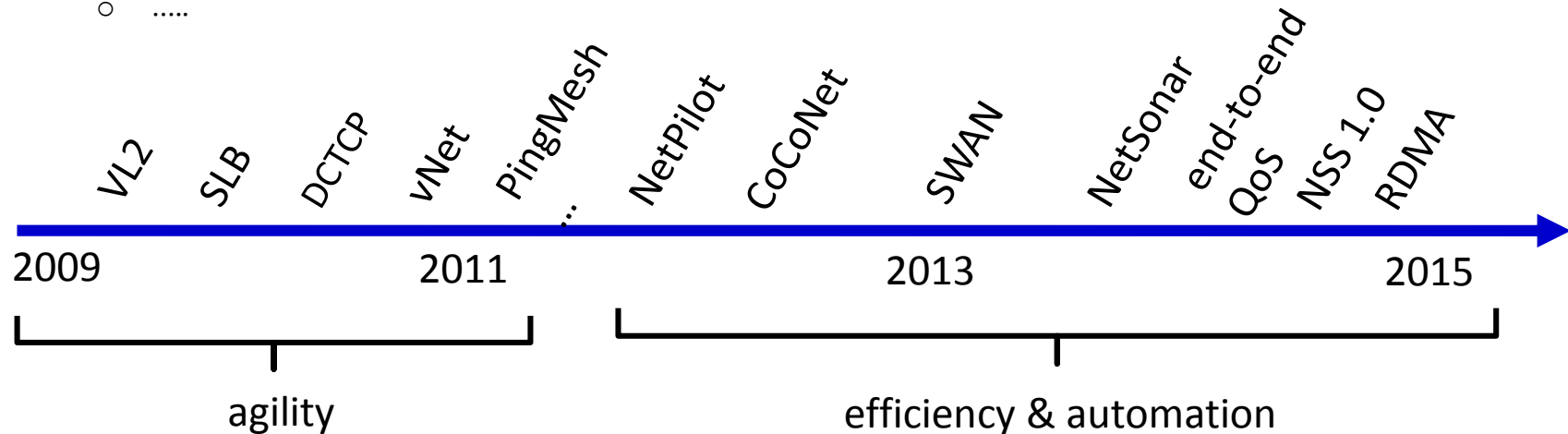
performance: significant number of circuits sit idle while others are oversubscribed (latency increases)

failures: long convergence time during network topology changes with planned and unplanned network events

MSR's contributions to Microsoft cloud networking & to academia

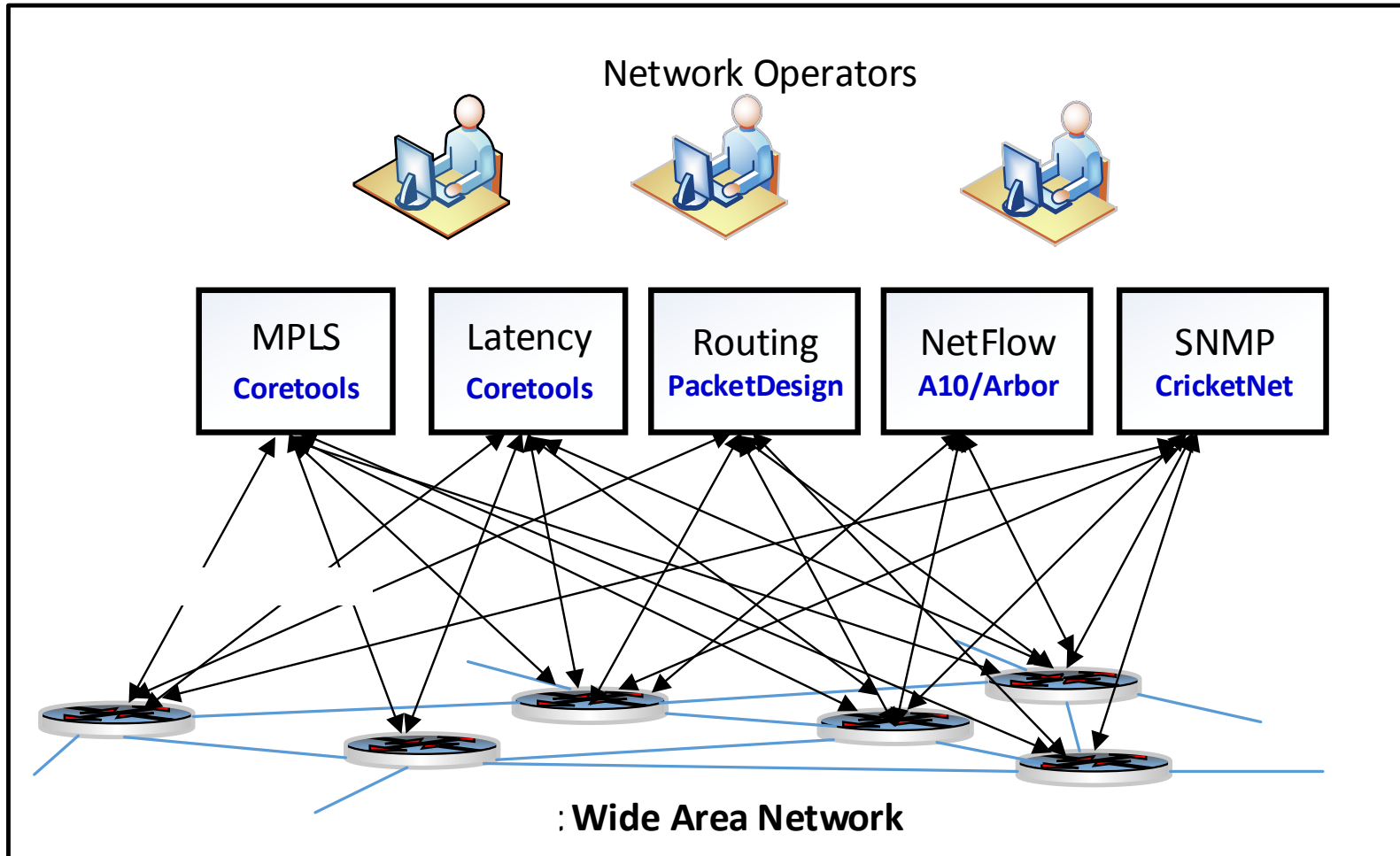
researchers worked hand-in-hand with Azure, Bing, Windows,

- steady stream of significant tech transfers
 - full-bisection bandwidth (Q10): 80x cost reduction, 20x outage reduction, in all Azure DCs
 - software load balancer (SLB): 15x cost reduction, carries all Azure traffic
 - software-defined WAN: increased inter-DC bw utilization from ~40% to ~95%,
 - virtual networking: enabled MSFT hybrid cloud offering via HyperV virtual network product
 -

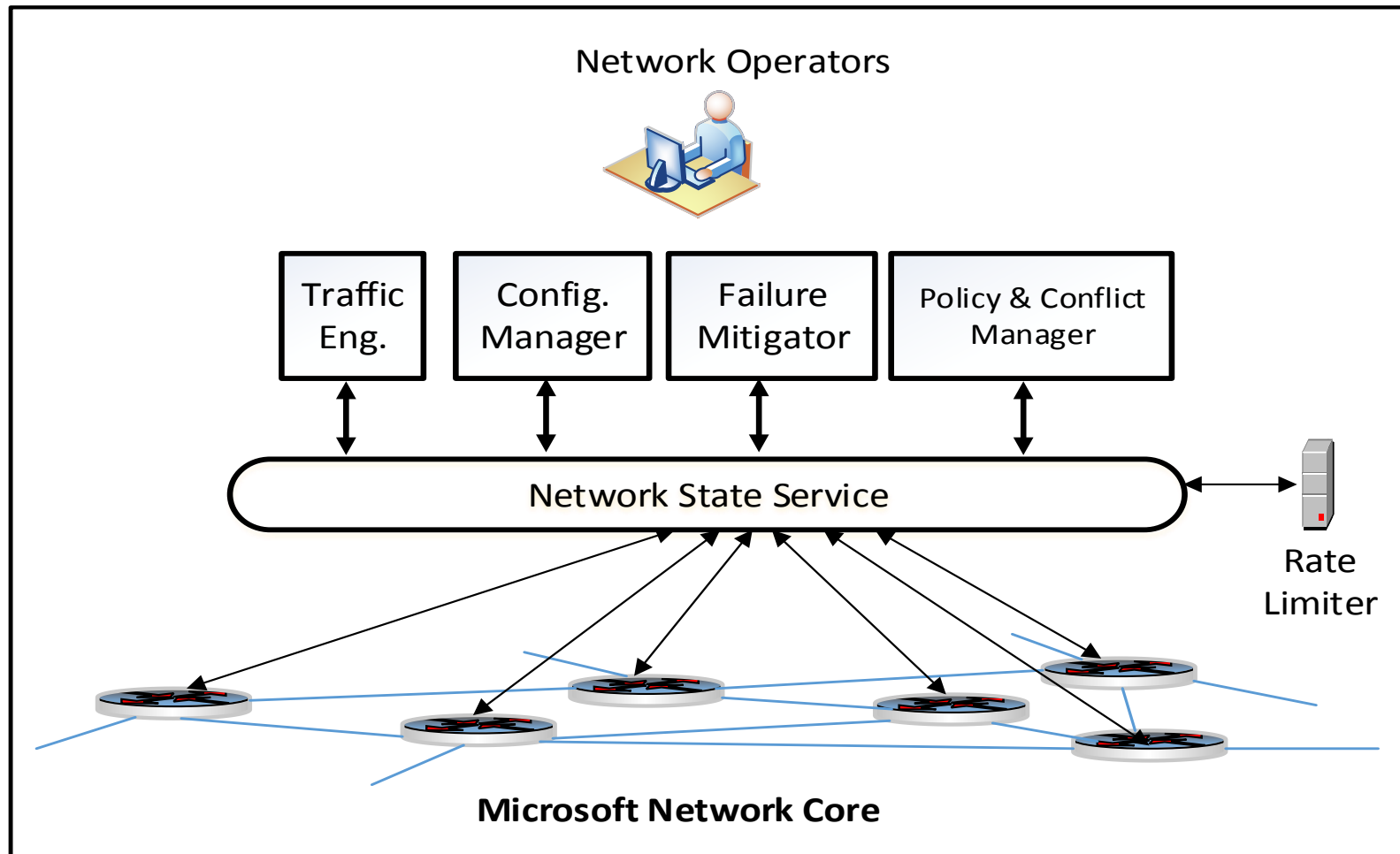


- plenty of research accolades as well
 - papers recognized as "Research Highlight" by ACM

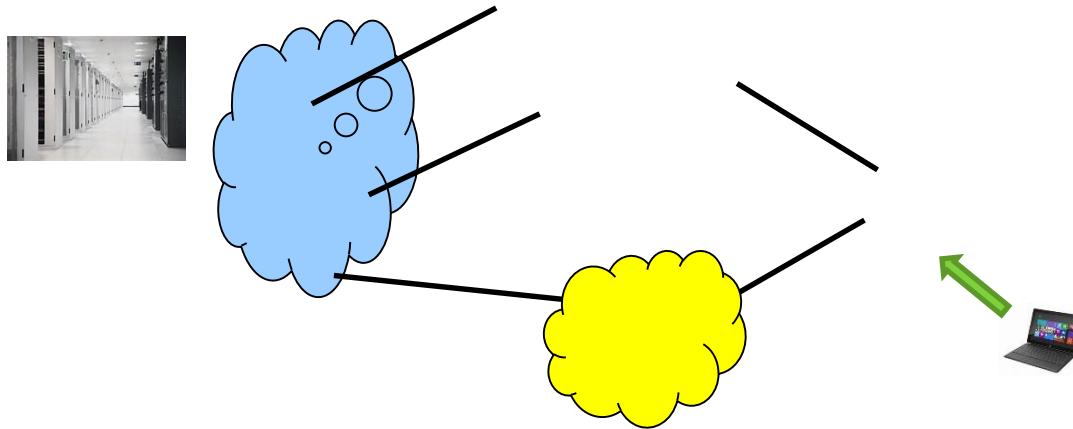
improving efficiency of wide area network



improving efficiency of wide area network with MSR's network state service



Internet: a network of networks of networks



as of March 6, 2013 (source: PEER 1)

- 22,961 AS numbers (AS numbers uniquely identify networks on the Internet, e.g. 8075 for Microsoft)
- 50,519 peering connections

... but we can reduce latency further

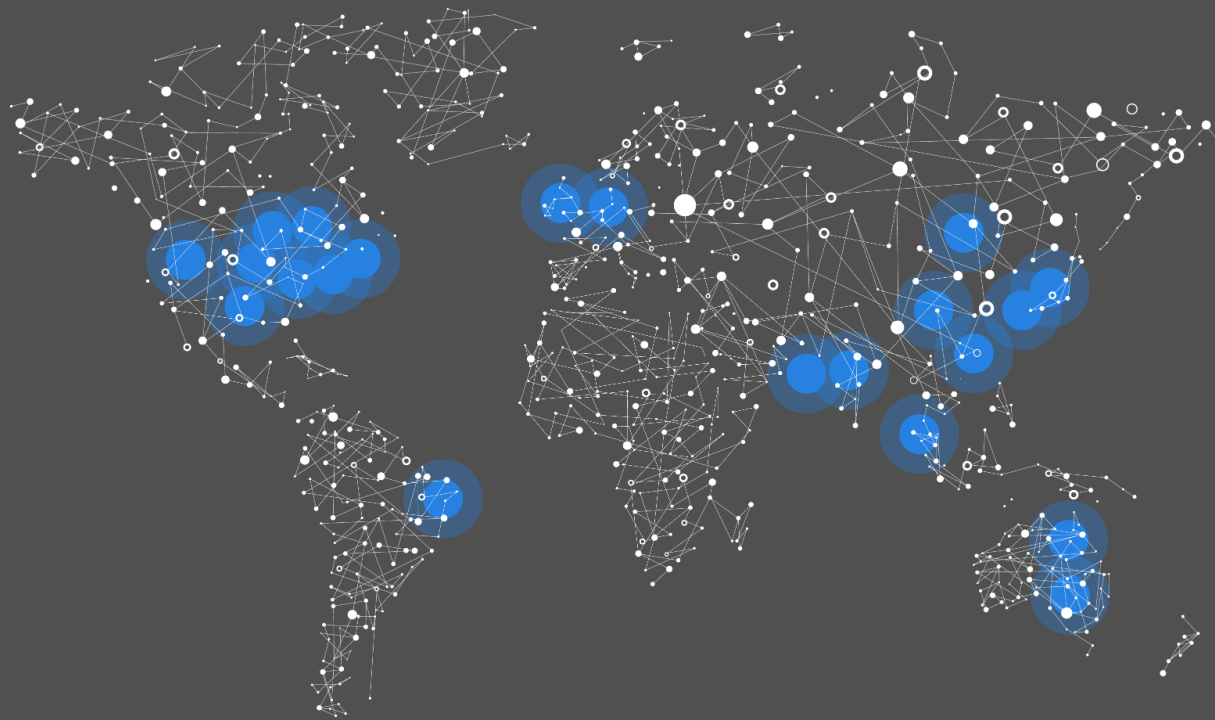
get the packets under our control as soon as possible

how?

- bring the cloud closer to the end-user
 - ✓ build lots of DCs around the world & place them in strategic locations

bringing the cloud closer

build lots of hyper-scale data centers around the world



 Azure compute regions

19

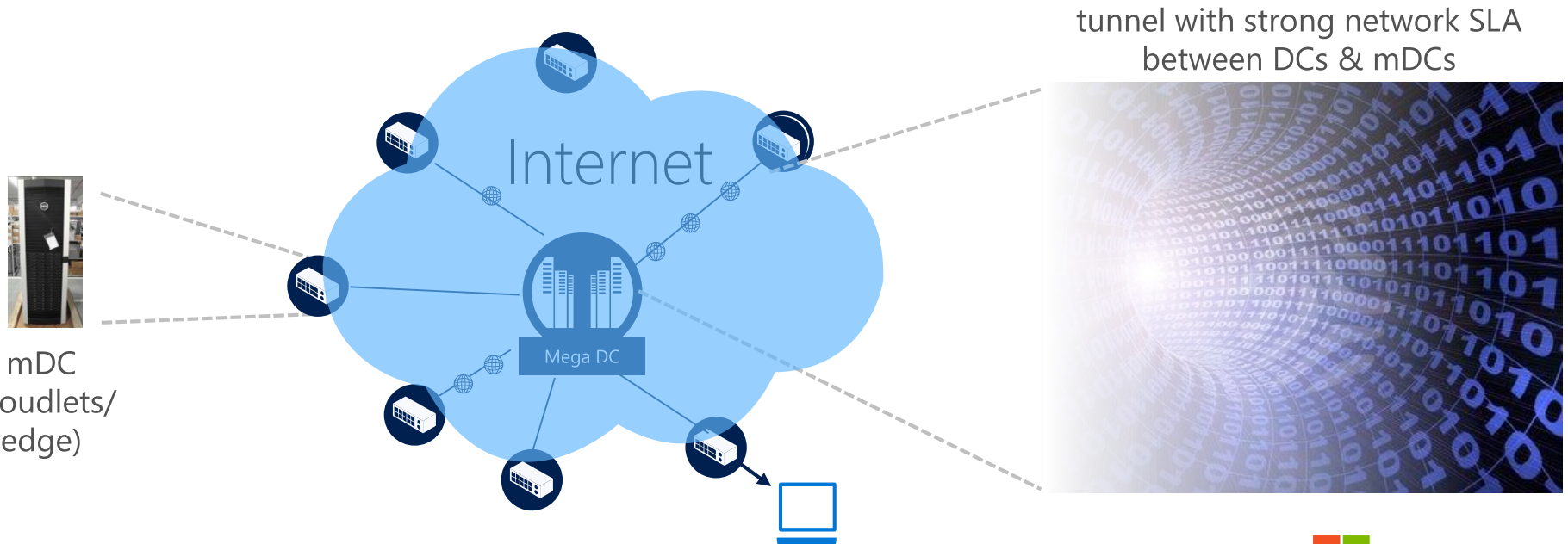
Azure compute
regions open today

more than AWS and
Google cloud combined

is building hyper-scale data centers enough?

no, it's capital intensive and expensive to operate

smarter approach: build an extensive infrastructure of micro DCs (1-10s of servers with several TBs of storage, \$20K-\$200K/mDC) and place them everywhere



micro DCs

site acceleration (classic)

content caching

- Xbox videos, NetFlix videos, Windows updates,...

split TCP connections

- from Bing data, on avg. can reduce latencies by ~30 msec
 - predictive search query responses improved ~25-35% based on random sampling before and after deploying edge serves in a couple of US cities

Akamai
Limelight
CloudFront
Level 3
EdgeCast
Rackspace
:
:

mDCs are “classic” CDNs nodes, that can improve the performance of search engines, office productivity tools, video and audio conferencing & future cloud services

additional benefits of mDCs

latency reduction

- ✓ serve static content immediately
- ✓ SSL termination / split TCP
- edge to DC protocol enhancements

bandwidth saving

- ✓ compression
- procrastination
- edge analytics

service & internet monitoring

reliable connectivity

- overlay networking
- path diversity

battery saving

- computation offloads
- client proxying

high-end game streaming

- lower device cost
- reduce developer fragmentation

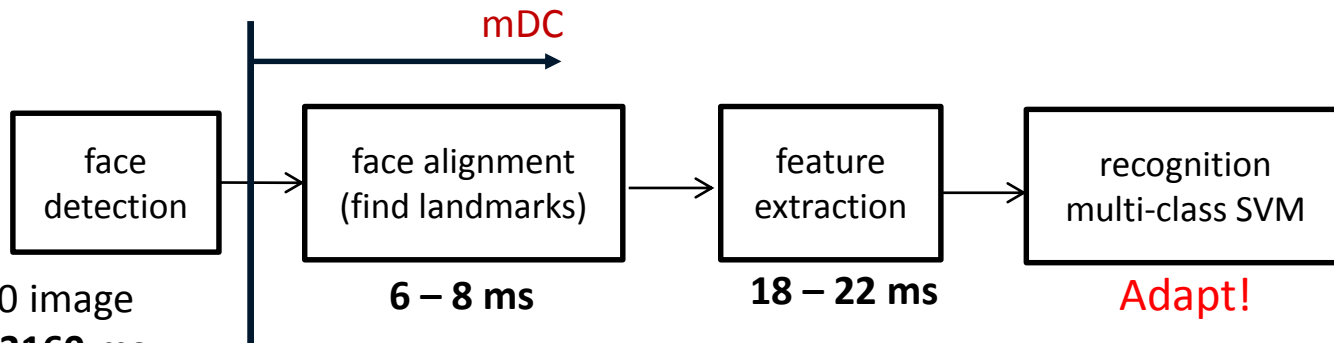
new services

protection against DoS

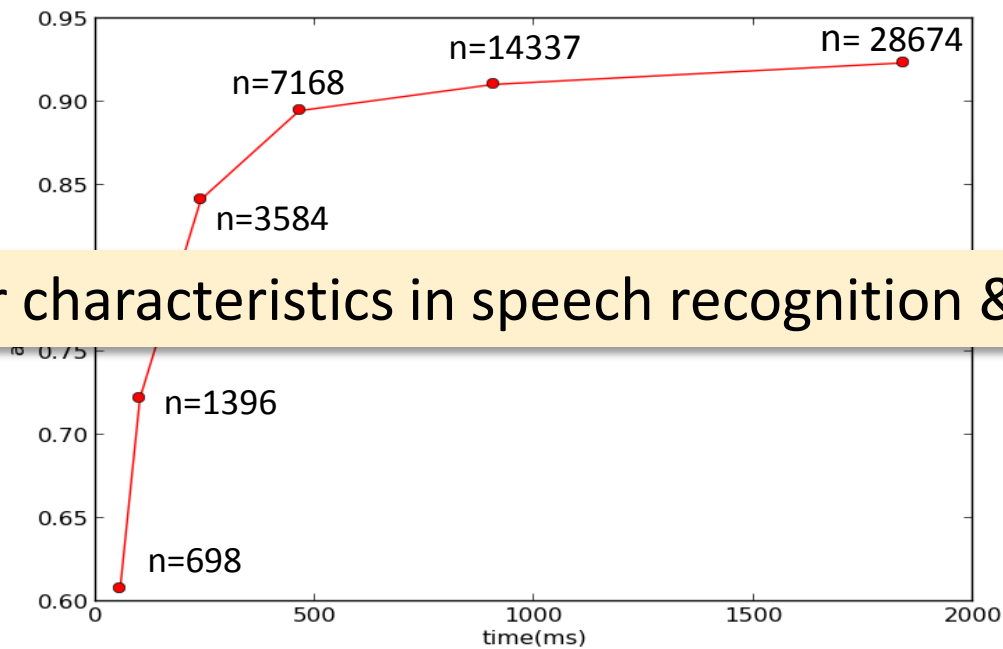
reduced load on DCs

new services: object recognition

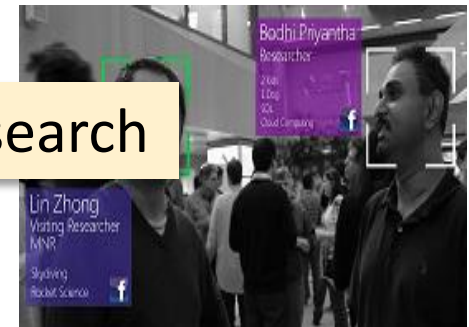
the lower the latency, the better the results



For a 640x480 image
client: 890 - 3160 ms
server: 72 - 115 ms



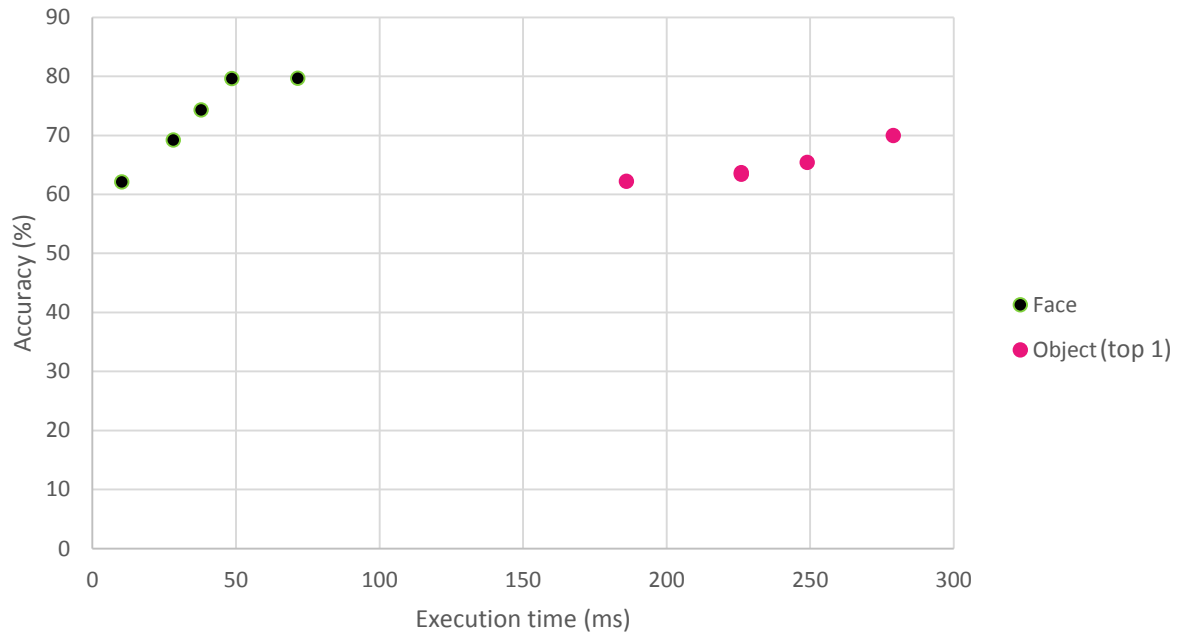
similar characteristics in speech recognition & search



Face prediction Time

using DNNs - similar results - lower transport latency helps

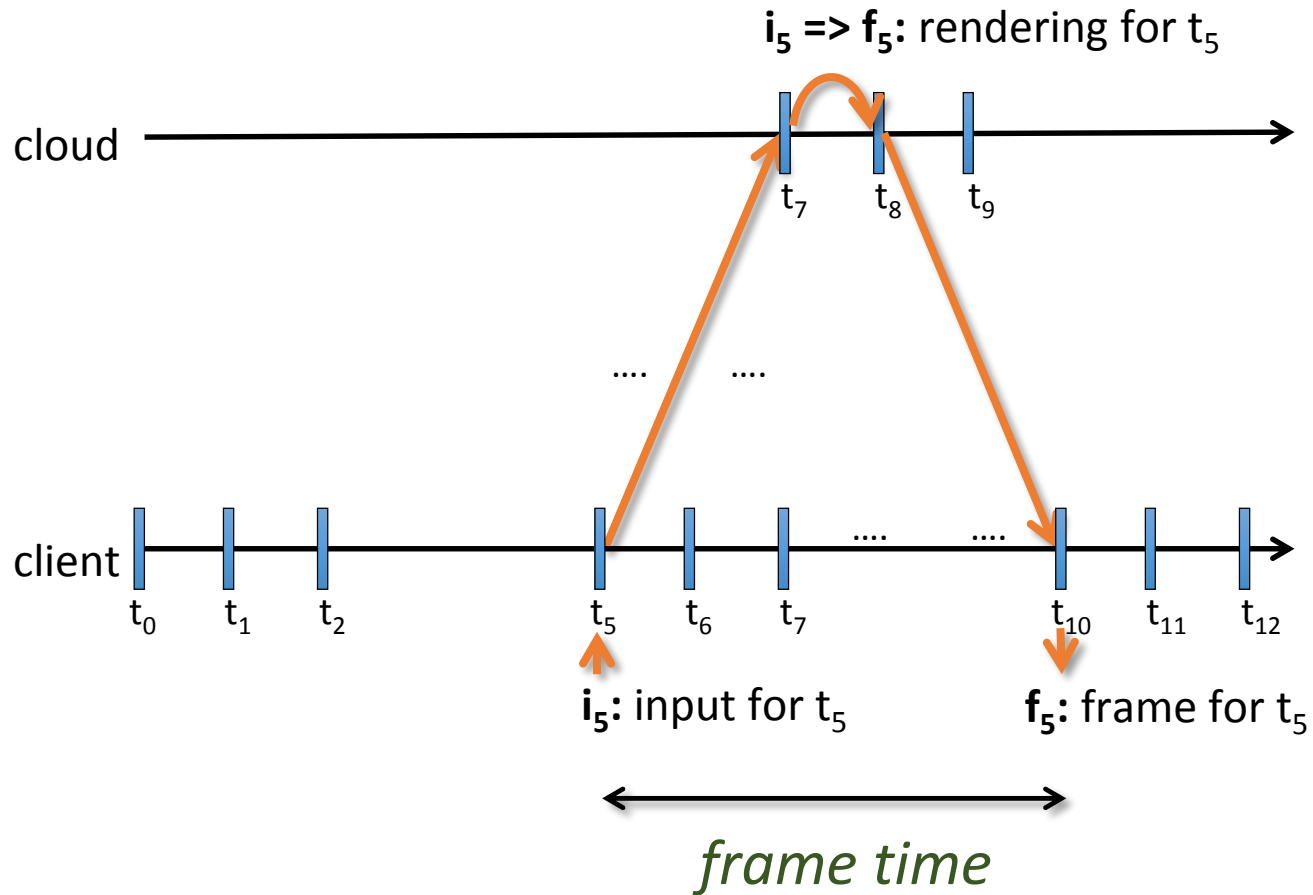
model execution time vs. accuracy (core i7)



50-100ms can allow ~10-20% more accurate model

face recognition with mDCs

(new) service: cloud gaming



cloud gaming

(with speculative execution)

Outatime

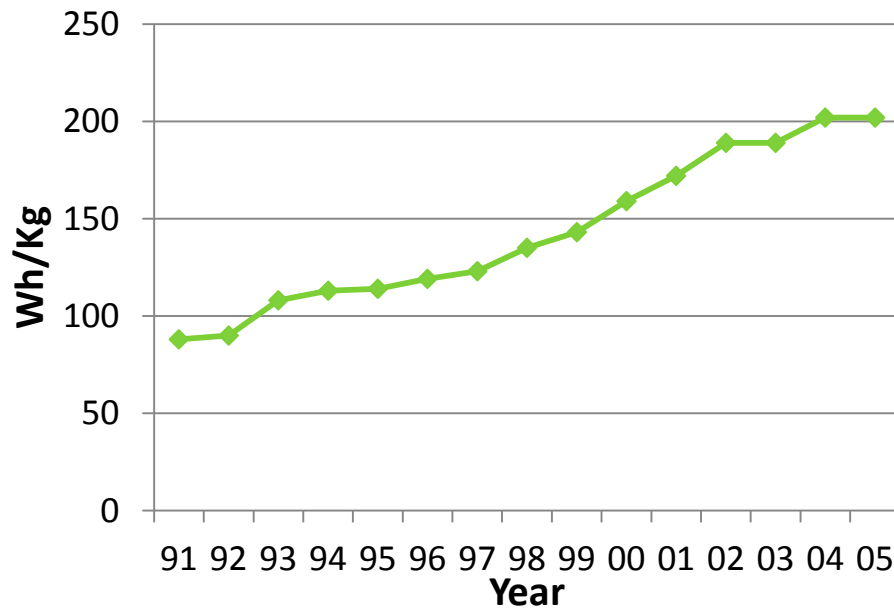
Improving Cloud Gaming w/
Speculative Execution



battery life...

silver bullet seems unlikely

Li-Ion energy density



lagged behind

- higher voltage batteries (4.35 V vs. 4.2V) – 8% improvement
- silicon anode adoption (vs. graphite) – 30% improvement

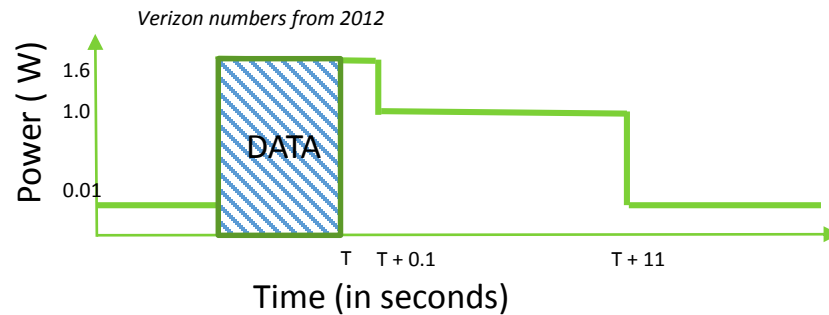
trade-offs

- fast charging = lower capacity
- slow charging = higher capacity

contrast with
CPU performance improvement during same period: 246x

battery use in SmartPhones...

LTE consumes $> 1.5W$ when active
LTE chip active for ~ 10 secs of extra tail time ($1W$ power)



....but how did we get here

a bit of context/history... 6 years ago

The New York Times

Customers Angered as iPhones Overload AT&T

By JENNA WORTHAM
Published: September 2, 2009

The New York Times

DIGITAL DOMAIN
AT&T Takes the Blame, Even for the iPhone's Faults

By RANDALL STROSS
Published: December 12, 2009

PCWorld

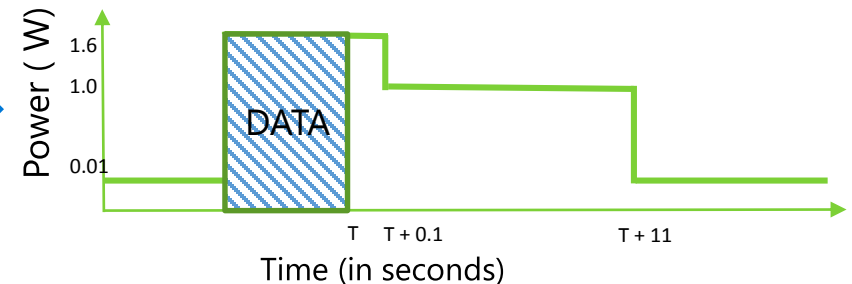
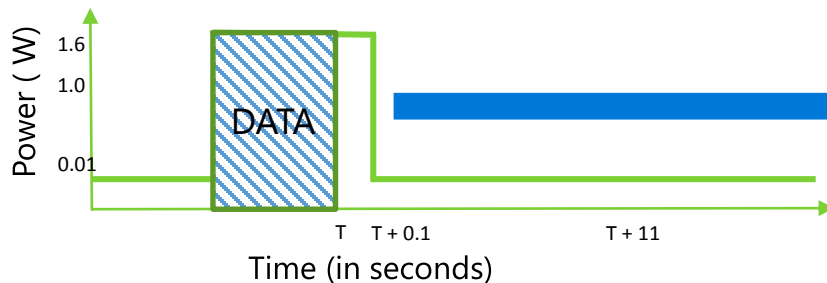
Report: AT&T Reputation Tarnished by iPhone Flaws

By Tony Bradley, PCWorld
Dec. 14, 2009 2:01 PM

original design:
bring radio to low power state immediately

mobile operator requirement:
keep LTE chip **active for ~10 sec.** of extra tail time
(to reduce the signaling load)

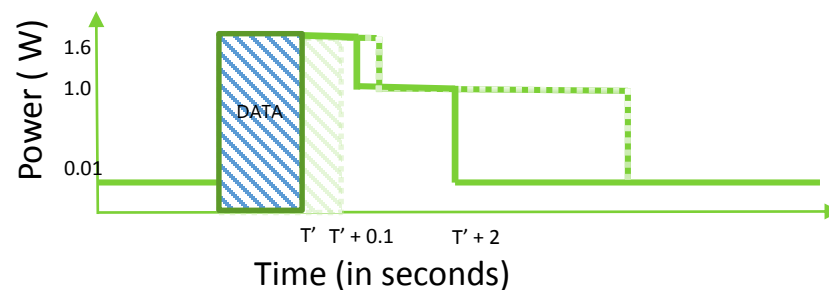
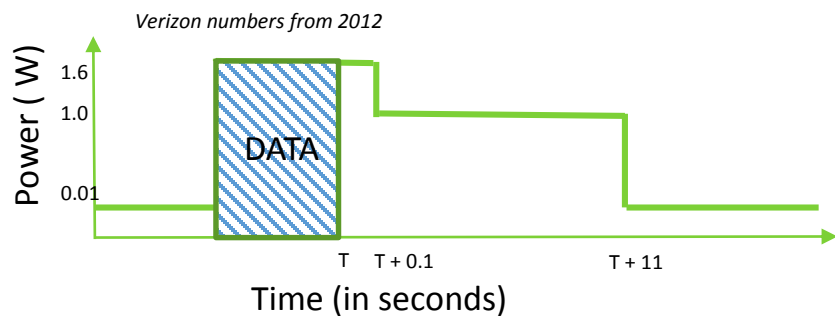
Verizon numbers from 2012



mDCs can increase use time

LTE consumes > 1.5W when active
LTE chip active for ~10 secs of extra tail
time (1W power)

with mDCs:
faster transfers => less time in high power state
aggressively enter lowest power state



Energy savings / transfer: $1.6W * \text{speedup} + 1W * 9\text{sec} = 10.6\text{J}$ (assuming speedup of 1 second)

for 20 network transfers/hour (notifications, email, etc.), with 1 sec speedup
total energy savings per 24 hr. = 6624 J

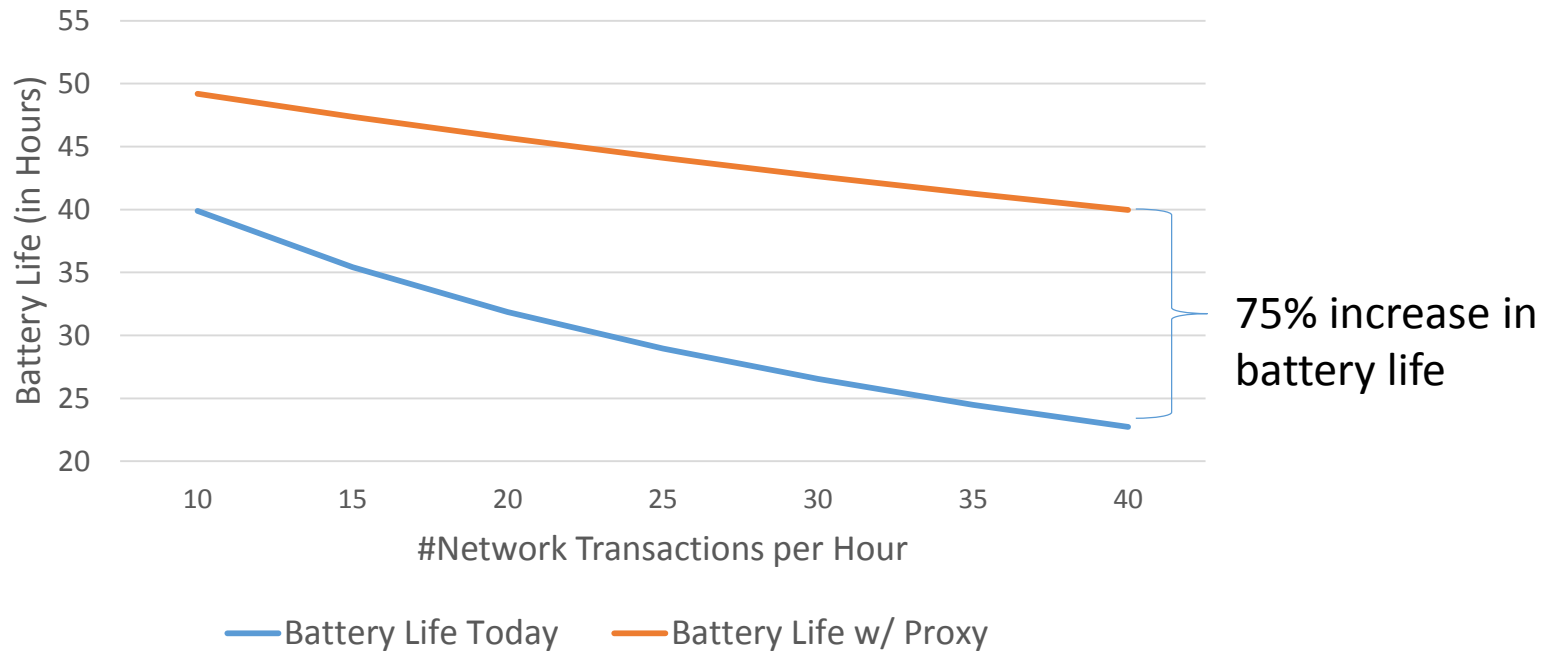
→ Saving of **26%** in a 1500 mAH cell phone battery*

* Samsung Standard LI-ION battery with rating of 1500mAh/3.7Vdc

especially good for mobile battery life improvement



calculated for a 30 msec speedup / network transaction



these types of saving occur across the board for all battery types and all types of mobile devices

* Samsung Standard LI-Ion battery with rating of 1500mAh/3.7Vdc

saving bandwidth....



security,
traffic,
tracking



locating
objects of
interest



customer
queue
analytics

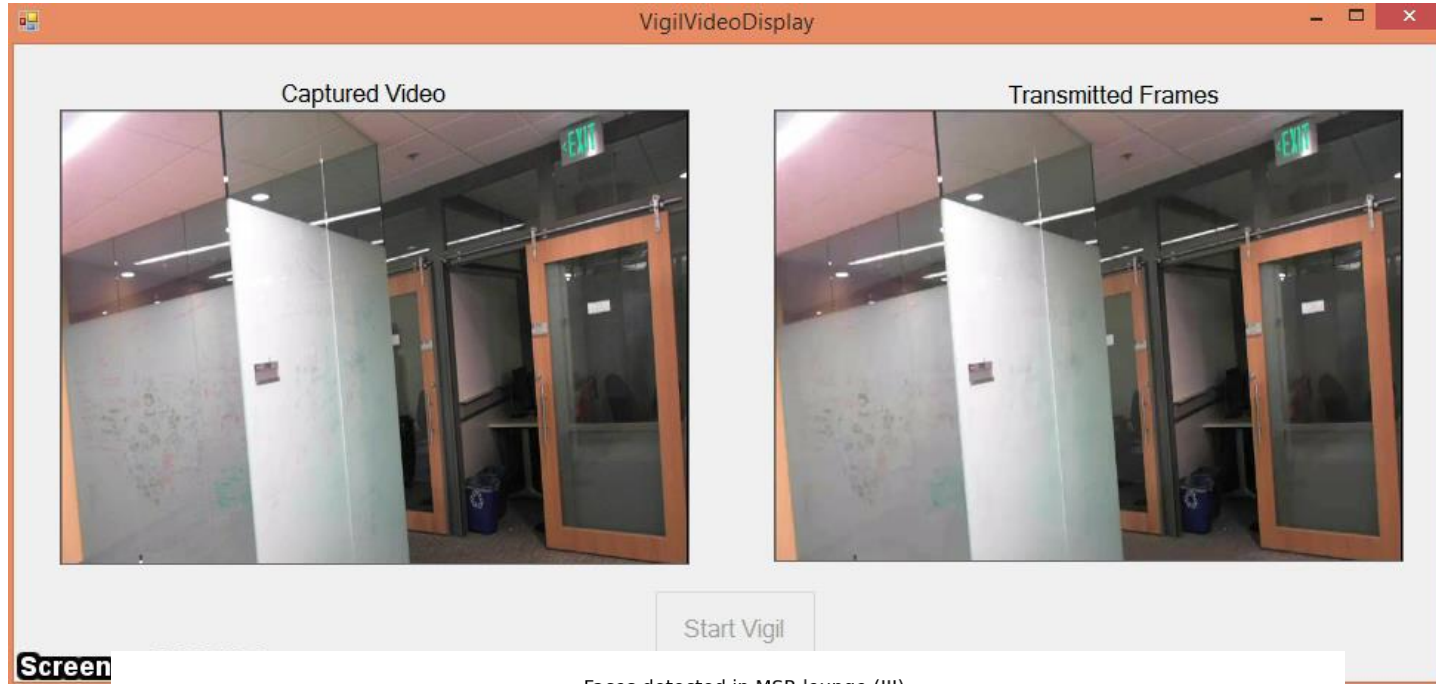
current approach

- upload the captured video to the cloud for remote analysis

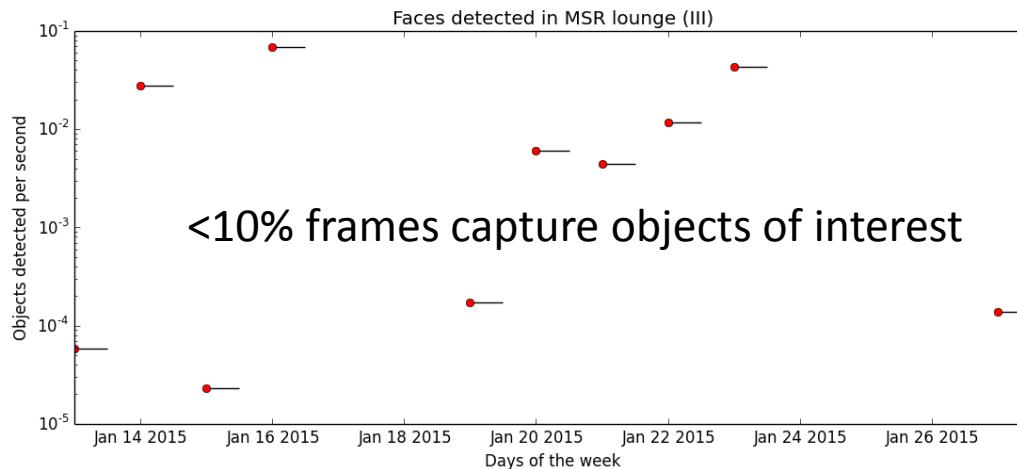
observations

- too much data captured per hour (>10GB/hour)
- bandwidth limits scale and use of system
- unable to support near real-time tracking & security

saving network bandwidth (wireless video surveillance)



Screen



saving network bandwidth (parking spot detector)



Microsoft
Research

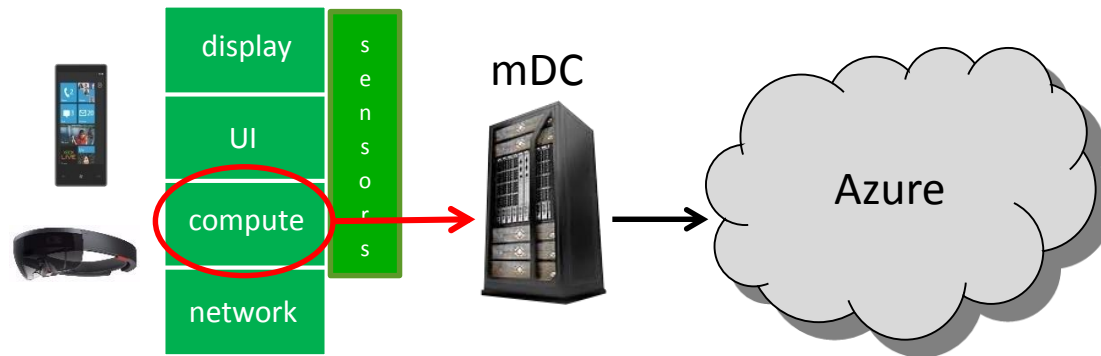
UPMC
SORBONNE UNIVERSITÉS



a couple of on-going problems

offloading computation

remote execution reduces energy consumption and improves performance



challenges

- what to offload?
- how to dynamically decide when to offload?
- **how to do so with minimum programmer effort?**
- **how to support multi-tenancy with bullet-proof privacy?**

programming frameworks for cloud offloading

| | Microsoft's MAUI | Intel's CloneCloud | USC's Odessa |
|-----------------------|------------------|--------------------|--------------|
| remote execution unit | methods | threads | tasks |

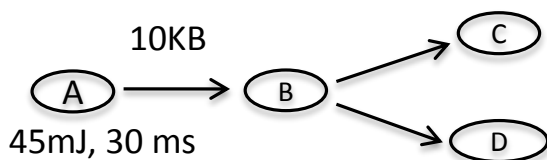
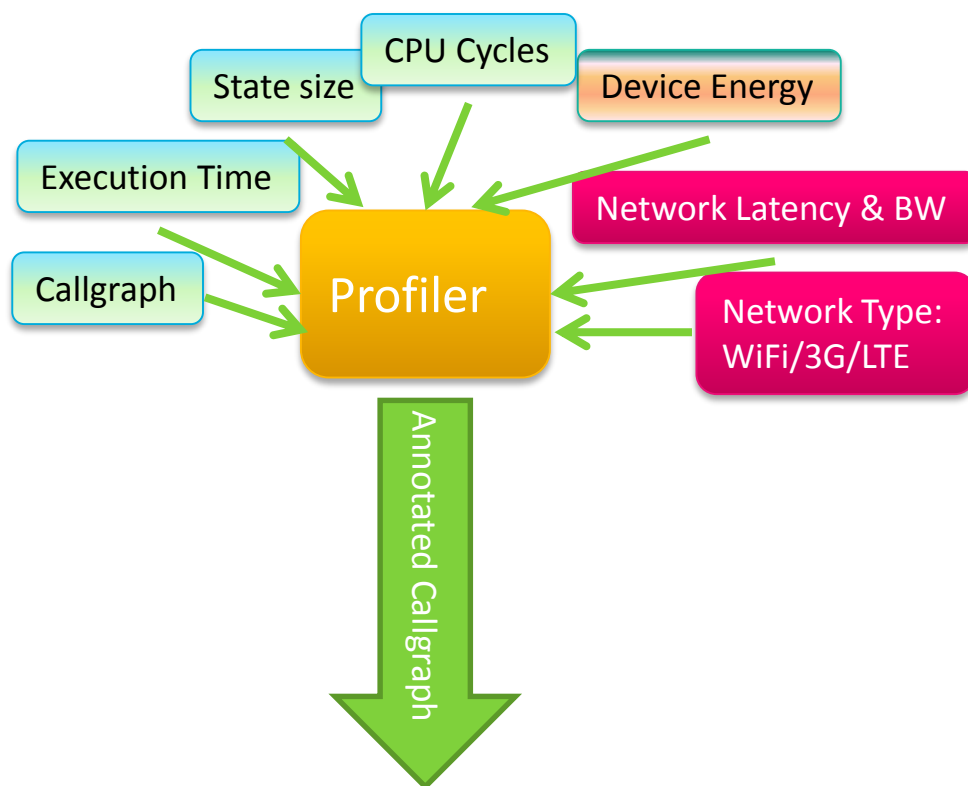
- MAUI exploits .NET framework to dynamically partitioning & offload method execution [MobiSys'10]
- CloneCloud supports existing applications, but requires tight synchronization between cloud and phone [EuroSys 2011]
- Odessa creates a data-flow graph to exploit parallelism [MobiSys 2011]

all have a **profiler** & a **solver**

also see: <http://elijah.cs.cmu.edu/>

MAUI's profiler and decision engine

profiler



decision engine:

partition a running app – use Integer Linear Program

Example – Maximize:

$$\sum_{v \in V} (I_v \times E_v) - \sum_{(u,v) \in E} (|I_u - I_v| \times C_{u,v})$$

energy saved cost of offload

Such that:

$$\sum_{v \in V} (I_v \times T_v) + \sum_{(u,v) \in E} (|I_u - I_v| \times B_{u,v}) \leq \text{Lat.}$$

execution time time to offload

and

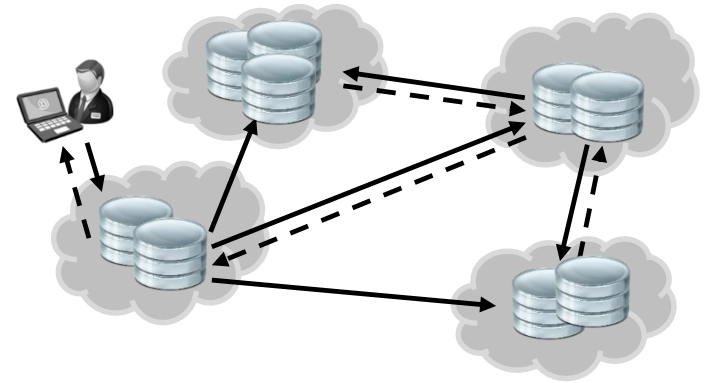
$$I_v \leq R_v \text{ for all } v \in V$$

geo-distributed analytics

lots of data being generated at the edges, need support for sophisticated analysis

possible solution(s)

- pull all data into a central data center; answer queries from there
- leave data where it is collected; fetch on demand per query



costly and wasteful; not realtime
very long latency; can't run Hive or Spark on WAN

Observations

- connectivity is expensive, low bw & high latency
- need to support near real-time triggers (e.g.. faults/ fire)
- some of the data is infrequently accessed

geo-distributed analytics

allow data & query tasks to be placed at any site

- some datasets remain at the edge; others move to resource-rich DCs
- make job schedulers' robust to high latency by pipelining

mimic *optimal* data & task placement

- minimize average query latency
 - E.g., move data iff the cumulative *shuffle volume* of its queries exceeds data size
 - Eg., place network-heavy tasks on a site where there is more data to be read

recapping benefits of mDCs

latency reduction

- serve static content immediately
- SSL termination / split TCP
- edge to DC protocol enhancements

bandwidth saving

- compression
- procrastination
- edge analytics

service & internet monitoring

reliable connectivity

- overlay networking
- path diversity

battery saving

- computation offloads
- client proxying

high-end game streaming

- lower device cost
- reduce developer fragmentation

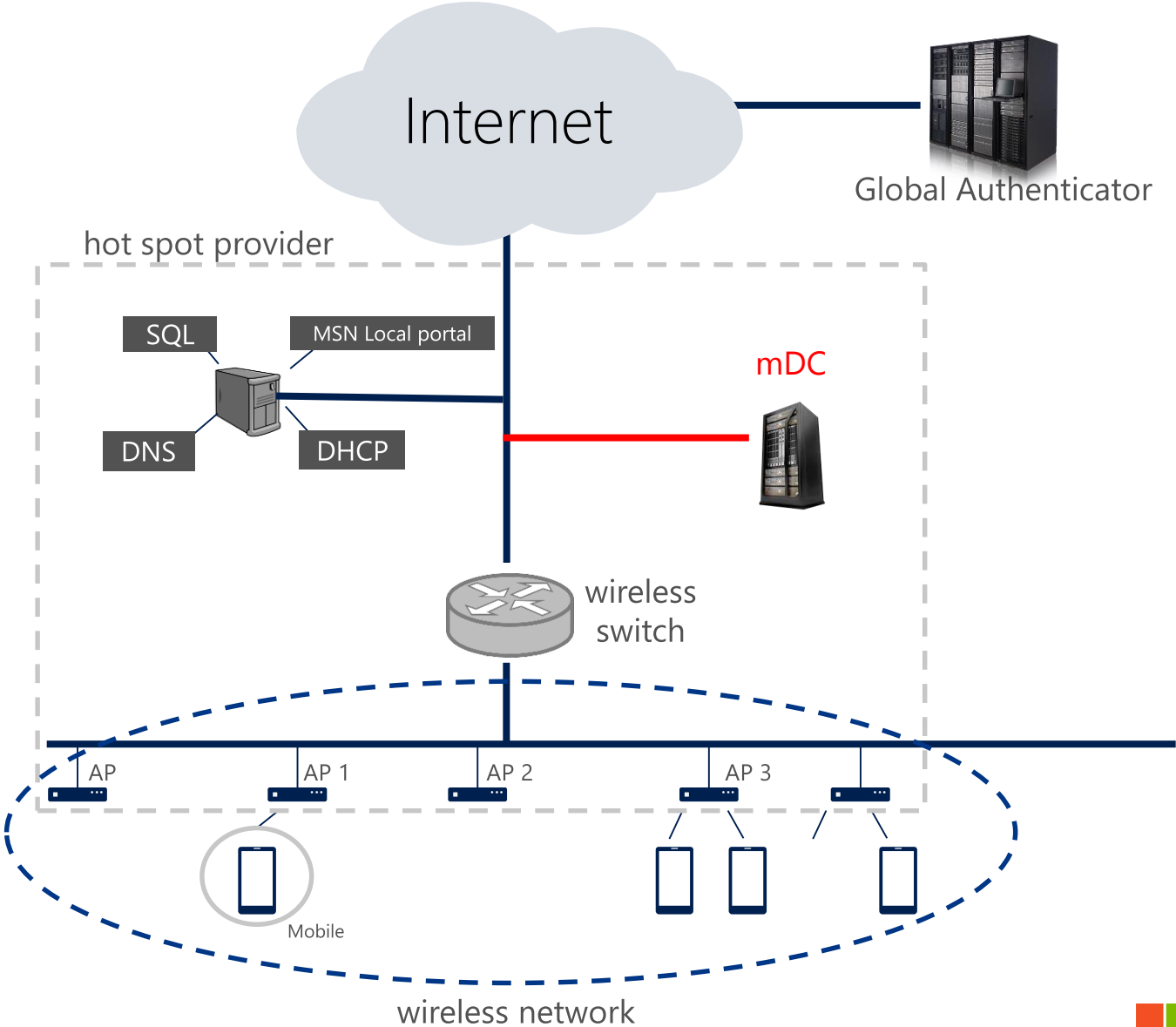
new services

protection against DoS

reduced load on DCs

deployment

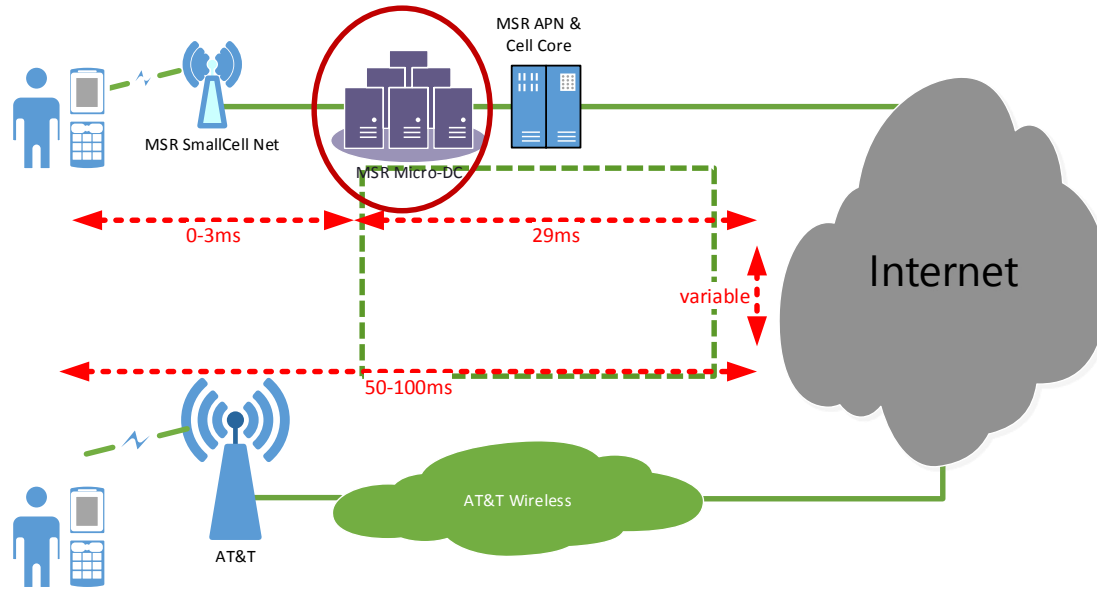
mDCs with Wi-Fi or White-Fi



mDC with small cells



QCOM's Small Cell



| | |
|----------|-----------|
| Downlink | ~110 Mbps |
| Uplink | ~15 Mbps |
| RTT | ~10 msec |

```

Telnet 127.0.0.1
C:\>tracert any.edge.bing.com
Tracing route to any.edge.bing.com [204.79.197.200]
over a maximum of 30 hops:
  1  37 ms  34 ms  39 ms  172.26.241.113
  2  *      37 ms  *      172.26.236.2
  3  38 ms  38 ms  43 ms  172.26.96.11
  4  38 ms  39 ms  39 ms  172.26.96.193
  5  50 ms  41 ms  40 ms  172.18.3.241
  6  44 ms  37 ms  60 ms  12.249.2.25
  7  44 ms  43 ms  44 ms  12.83.180.6
  8  48 ms  47 ms  42 ms  12.83.180.14
  9  45 ms  52 ms  44 ms  cr81.st0wa.ip.att.net [12.122.5.197]
 10  93 ms  120 ms  43 ms  12.122.111.9
 11  45 ms  44 ms  46 ms  12.249.36.6
 12  *      *      *      Request timed out.
 13  *      *      *      Request timed out.
 14  *      *      *      Request timed out.
 15  50 ms  50 ms  50 ms  origin.any.bing.com [204.79.197.200]
Trace complete.
C:\>
    
```

tracert from AT&T LTE to any.edge.bing.com (15 hops)

```

Command Prompt
C:\Users\sagarwal>tracert any.edge.bing.com
Tracing route to any.edge.bing.com [204.79.197.200]
over a maximum of 30 hops:
  1  *      *      *      Request timed out.
  2  42 ms  27 ms  39 ms  131.107.151.1
  3  43 ms  98 ms  26 ms  ge-3-0-0-401.icar-sttlwa01-02.infra.pnw-gigapop.net [209.124.190.238]
  4  35 ms  27 ms  39 ms  ae1-706.iccr-sttlwa01-03.infra.pnw-gigapop.net [207.231.240.1]
  5  32 ms  28 ms  27 ms  microsoft-1-lo-jmb-706.sttlwa.pacificwave.net [207.231.240.7]
  6  30 ms  29 ms  27 ms  ae0-0.wst-96che-1a.ntwk.msn.net [204.152.140.105]
  7  *      *      *      Request timed out.
  8  *      *      *      Request timed out.
  9  *      *      *      Request timed out.
 10  43 ms  27 ms  38 ms  any.edge.bing.com [204.79.197.200]
Trace complete.
C:\Users\sagarwal>
    
```

tracert from SC to any.edge.bing.com (10 hops)



the wave is coming ...



Dynamically manipulate images in the Cloud for responsive web design

Fast and Easy Access Control for your Web Applications



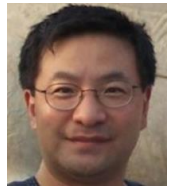
Increasing Mobile Operators' Value Proposition With Edge Computing

Turn bit pipes into smart pipes with an Intel® architecture-based server embedded into a Nokia Siemens Networks* base station

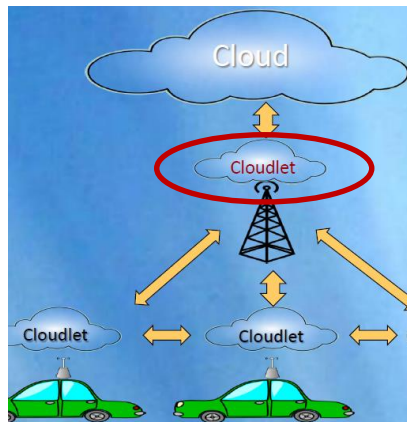


“local cloud are essential for backbone and core network scalability”

Dr. Geng Wu, Chief Scientist, **Intel** (Wireless World Research Forum, Vancouver, BC, Oct. 22, 2013)



5G with **Undelay Networks** and **Local Cloud**



“cloudlets for reducing latency, security and reliability”

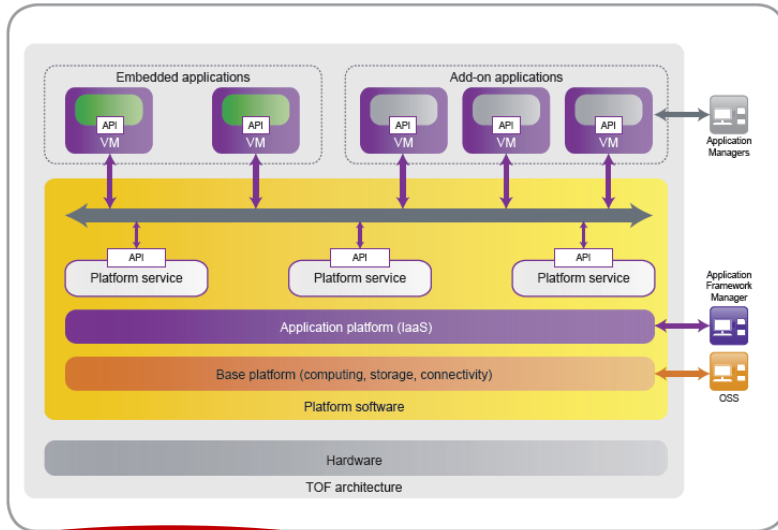
- Dr. David Soldani, VP **Huawei**
- (IEEE ICC, June 12, 2013)



...and it's becoming bigger

MOs moving towards edge services

Liquid Net



News
Nokia Siemens to merge cloud, base-station computing to boost performance
 The company's Liquid Applications platform will use computing power in the cloud and in base stations, based on conditions.
 By Stephen Lawson, IDG News Service
 February 24, 2013 04:06 PM ET

Add a comment Print Share +1 Like 0 More

IDG News Service - Nokia Siemens Networks will expand the role of cellular base stations with a new platform that will store and deliver some application data locally, while tapping into information about subscribers and traffic to improve the process.

The company announced the system, called Liquid Applications, at an event in Barcelona on the eve of Mobile World Congress. Liquid Applications can improve consumers' mobile experience but cutting delays as well as delivering more relevant content, CEO Rajeev Suri

Figure 4: Base station application architecture



Nokia Networks reveals ETSI mobile edge computing collaboration

IBM, Intel, Vodafone, and Huawei all on board

October 20, 2014 | By Michael Carroll

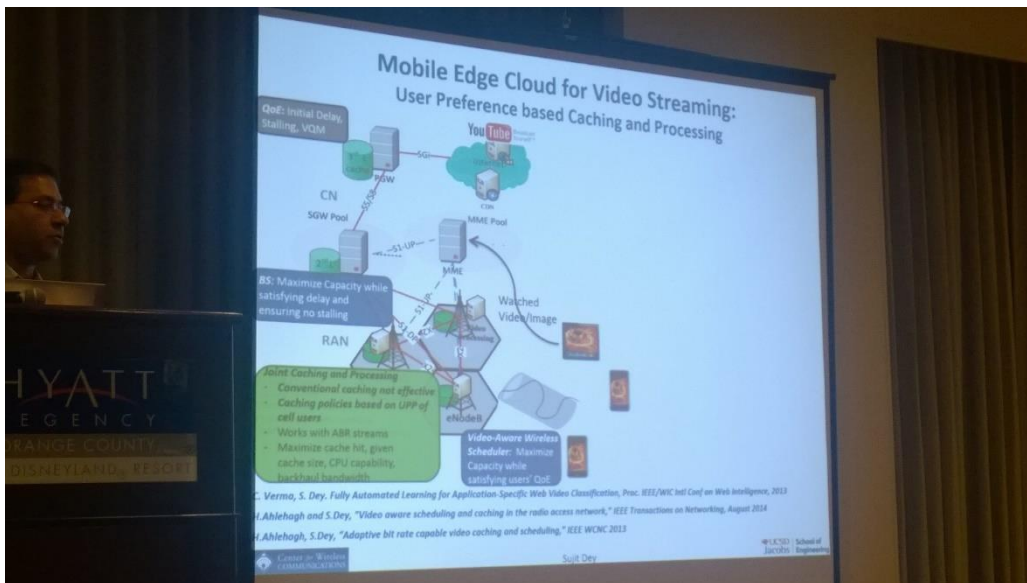
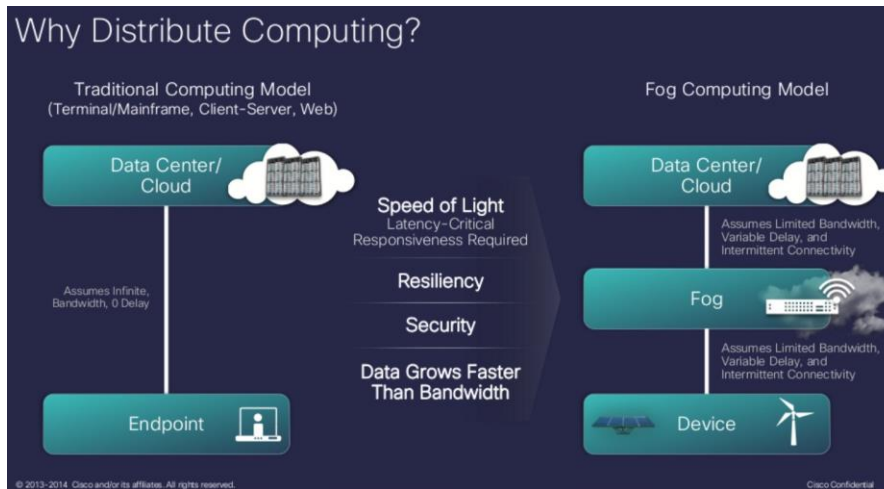


overheard at a recent conference (IEEE ICNC 2015)

“fog computing”



John Apostolopoulos
CTO & VP, Cisco, USA



Sujit Dey, Professor/ Director
Center for Wireless Communications UCSD

it's hot in the research community as well...

there is plenty of research literature (incl. MSR's) that shows edge computing significantly enhances mobile experience

first paper → Satya (CMU), Bahl (Microsoft), Caceres (AT&T), Davies (Lancaster)
The Case for VM-based Cloudlets in Mobile Computing
IEEE Pervasive Computing, October 2009

~ 900 citations

Cuervo (Duke), Balasubramanian (UMASS), Wolman, Saroiu, Chandra, Bahl (Microsoft)
MAUI: making smartphones last longer with code offload
ACM MobiSys conference, June 2010

~ 825 citations

[Why a Cloudlet Beats the Cloud for Mobile Apps](#)

Posted on December 13, 2009 by lewisshepherd

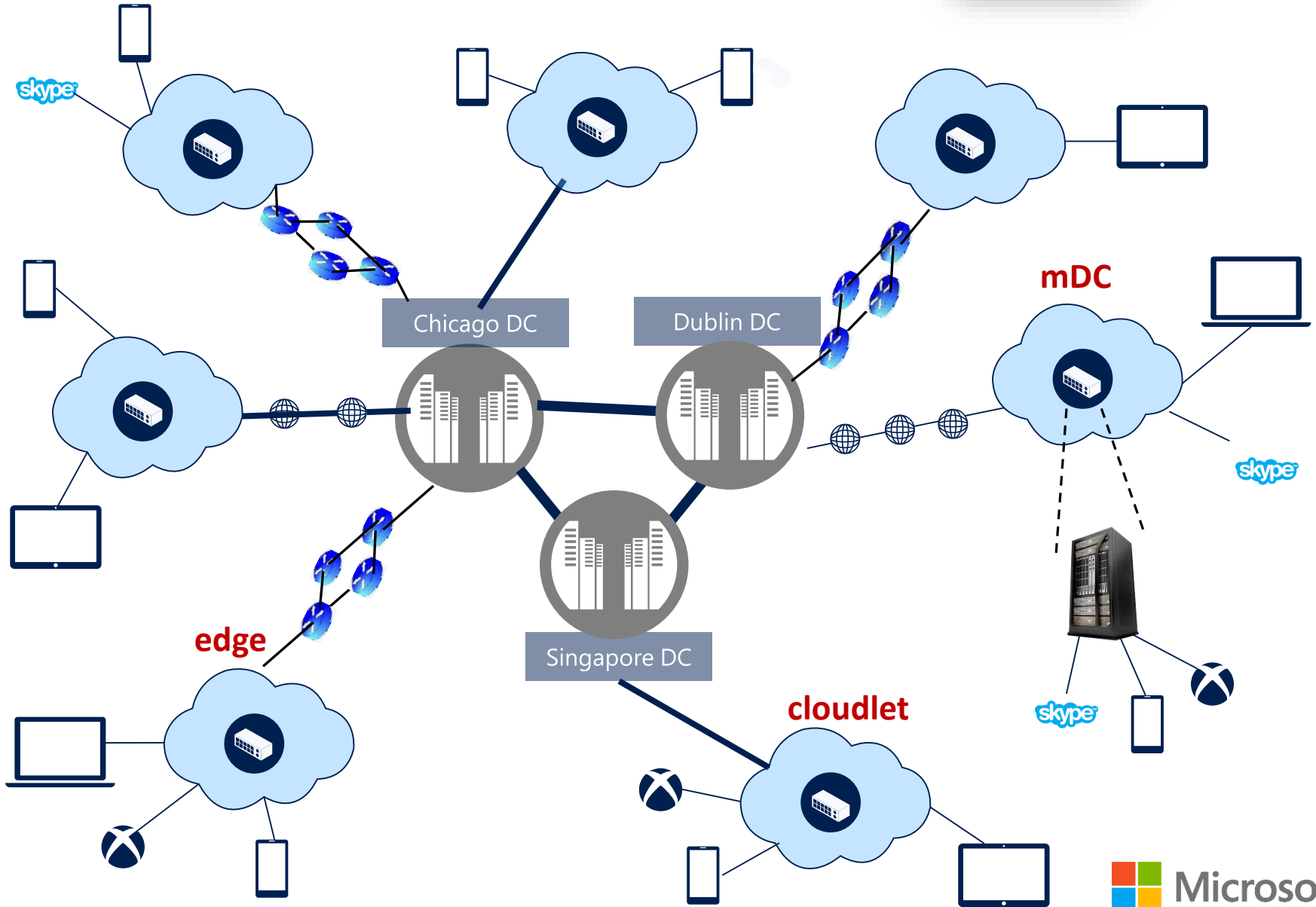
cloud computing 2020



=



hyper scale data center



with mDCs (cloudlets) you can...

- develop new (*latency sensitive, CPU & battery intensive*) (IoT) applications, which (*dynamically*) partition themselves
- pursue infrastructure research in an emerging cloud platform, which promises to be pervasive
- deploy your own mDCs & connect them to Azure

merci!



mDC benefits - app & game streaming

run any ecosystem's apps on resourced-starved devices by streaming them from the cloud

- circumvent client-side compatibility complexities
- with mDCs, reduce
 - latency -- keeping users engaged
 - jitter & packet loss – reduce user frustrating in highly interactive sessions
 - backbone bandwidth so both MOs and we pay less to other ISPs



note: standard proxy + split TCP insufficient for interactive traffic

mDCs can reduce dependency on cellular networks

offload to Wi-Fi aggressively

compress aggressively

} already doing this

procrastinate instead of prefetch

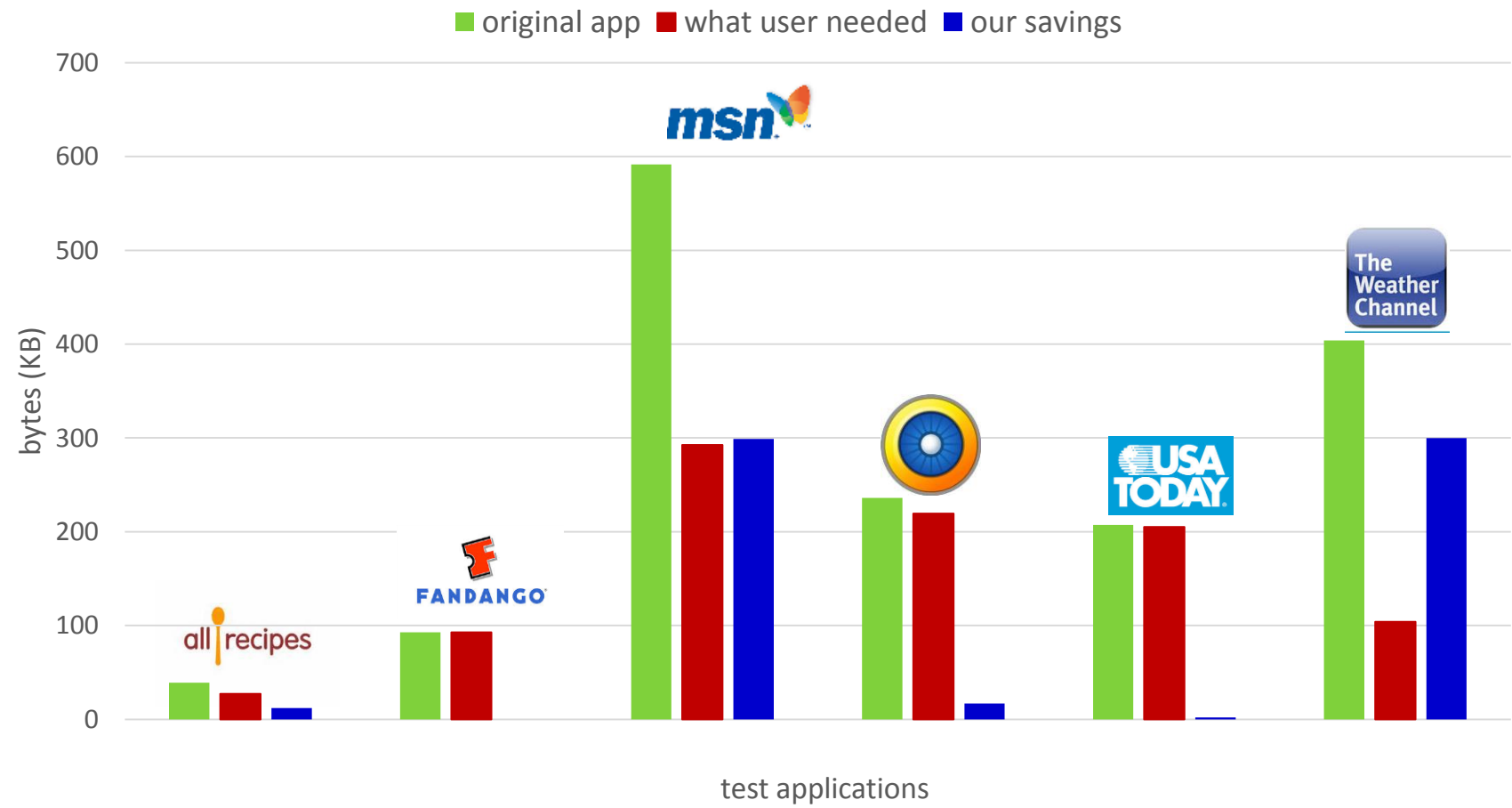
MobiSys 2014

- many network apps. fetch data whether or not it is consumed
- **idea:** mDC fetches the data but holds on to it until user explicitly needs it
 - ✓ save cellular bandwidth without the latency penalty

procrastinate & save

few results on bandwidth saving

the system automatically decides what is not needed by the end-user



micro datacenter - benefits

reducing dependency on cellular networks (with procrastination)

get data only when needed (**without mDC**)



get data only when needed (**with mDC**)

