# Advanced service-based data analytics: Models, Elasticity and APIs

## Hong-Linh Truong
## Distributed Systems Group, TU Wien

truong@dsg.tuwien.ac.at
dsg.tuwien.ac.at/staff/truong
@linhsolar

1

DISTRIBUTED SYSTEMS GROUP

# **Outline**

- Principles of elasticity for advanced service-based data analytics

- Data analytics within a single system

- Data analytics across multiple systems

- APIs management

DISTRIBUTED SYSTEMS GROUP

# Recall

data-as-a-service and data marketplaces are key elements for data-driven economy

# PRINCIPLES OF ELASTICITY FOR DATA ANALYTICS

DISTRIBUTED SYSTEMS GROUP

# Complex dependencies in (big) data analytics



- **More data** → more computational resources (e.g. more VMs)

- **More types of data** → more computational models → more analytics processes

- Change **quality of analytics**
  - Change quality of data
  - Change response time
  - Change cost
  - Change types of result (form of the data output, e.g. tree, visual, story, etc.)

DISTRIBUTED SYSTEMS GROUP

# Complex dependencies in (big) data analytics



a) Data analytics with the same type of data

b) Data analytics with multiple types of data

c) Impact of quality of results on data, analytics process, and computational models

ASE Summer 2016

Hong-Linh Truong, Schahram Dustdar, "Principles of Software-defined Elastic Systems for Big Data Analytics", (c) IEEE Computer Society, IEEE International Workshop on Software Defined Systems, 2014 IEEE International Conference on Cloud Engineering (IC2E 2014), Boston, Massachusetts, USA, 10-14 March 2014

DISTRIBUTED SYSTEMS GROUP

Elasticity principles can be used to support dynamic quality of analytics

DISTRIBUTED SYSTEMS GROUP

# Elasticity Principles: Elasticity of data and computational models

- Multiple types of objects from different sources with complex dependencies, relevancies, and quality

- Different data and computational models for the same analytics subject

- New analytics subjects can be defined and analytics goals can be changed

- Decide/select/define/compose not only computational models for analytics subjects but also data models based on existing ones

Management and modeling of elasticity of data and computational model during the analytics

DISTRIBUTED SYSTEMS GROUP

# Elasticity Principles: Elasticity of data resources

- Data provided, managed and shared by different providers

- Data associated with different concerns (cost, quality of data, privacy, contract, etc.

- Static data, open data, data-as-a-service, opportunistic data (from sensors and human sensing)

- Not just centralized big data and total data ownership

Data resources can be taken into account in an elastic manger: similar to VMs, based on their quality, relevancy, pricing, etc.

DISTRIBUTED SYSTEMS GROUP

# Elasticity Principles: <span style="color:red">Elasticity of humans and software as computing units</span>

- Human in the loop to solve analytics tasks that software cannot solve

- Human-based compute units can be scaled up/down with different cost, availability, and performance models

- Human-based compute units + software-based compute units for executing computational models

- Elasticity controls can be also done by humans

Provisioning hybrid compute units in an elastic way for computational/data/network tasks as well as for monitoring/control tasks in the analytics process

# Elasticity Principles: Elasticity of quality of analytics

- Definition of quality of analytics
  - Trade-offs of time, cost, quality of data, forms of output

- Using quality of analytics to select suitable computational models, data resources, computing units

- Multi-level control for the elasticity based on quality of analytics

Able to cope with changes in quality of data, performance, cost and types of results at runtime
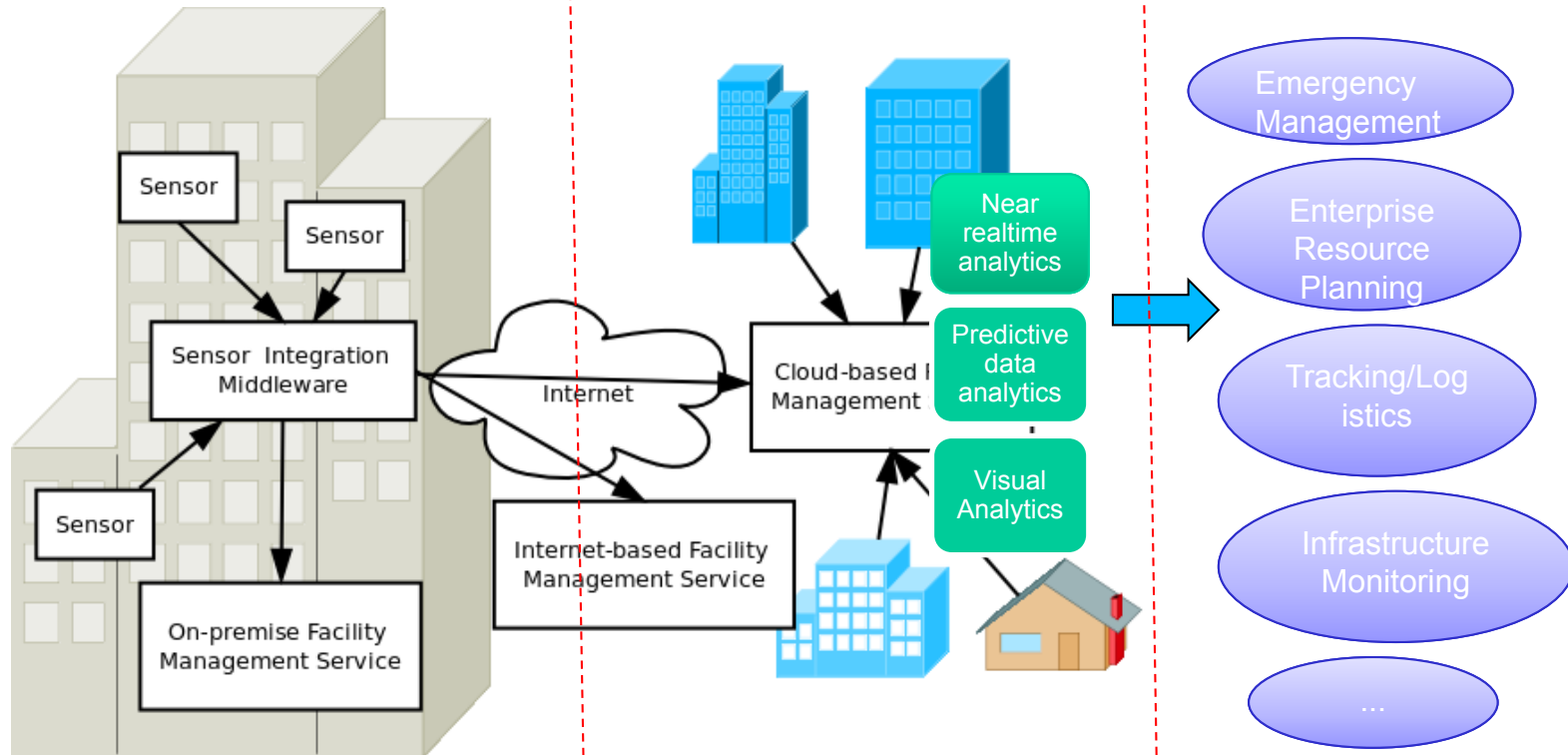
DISTRIBUTED SYSTEMS GROUP

# Advanced service-based analytics – which are fundamental engineering questions?

# Advanced service-based data analytics (1)



Infrastructure/Internet of Things

Internet/public cloud boundary

Organization-specific boundary

Sensor

Sensor

Sensor Integration Middleware

Sensor

On-premise Facility Management Service

Internet

Internet-based Facility Management Service

Cloud-based Management

Near realtime analytics

Predictive data analytics

Visual Analytics

Emergency Management

Enterprise Resource Planning

Tracking/Log istics

Infrastructure Monitoring

...

Cities, e.g. including:
10000+ buildings
1000000+ sensors

DISTRIBUTED SYSTEMS GROUP

# Advanced service-based data analytics -- fundamental concepts

Domain 1      Domain 2      Domain n

**Applications**

| Part A | → | Part B | → | ... | → | Part N |

**System infrastructures**

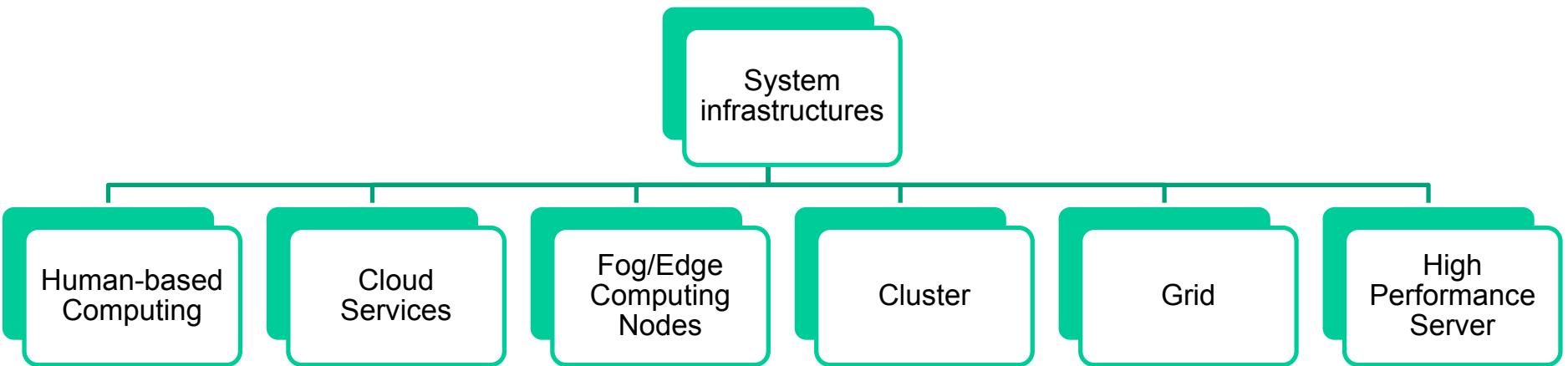IoT     Edge servers     Local Cloud     Public cloud

# Design questions

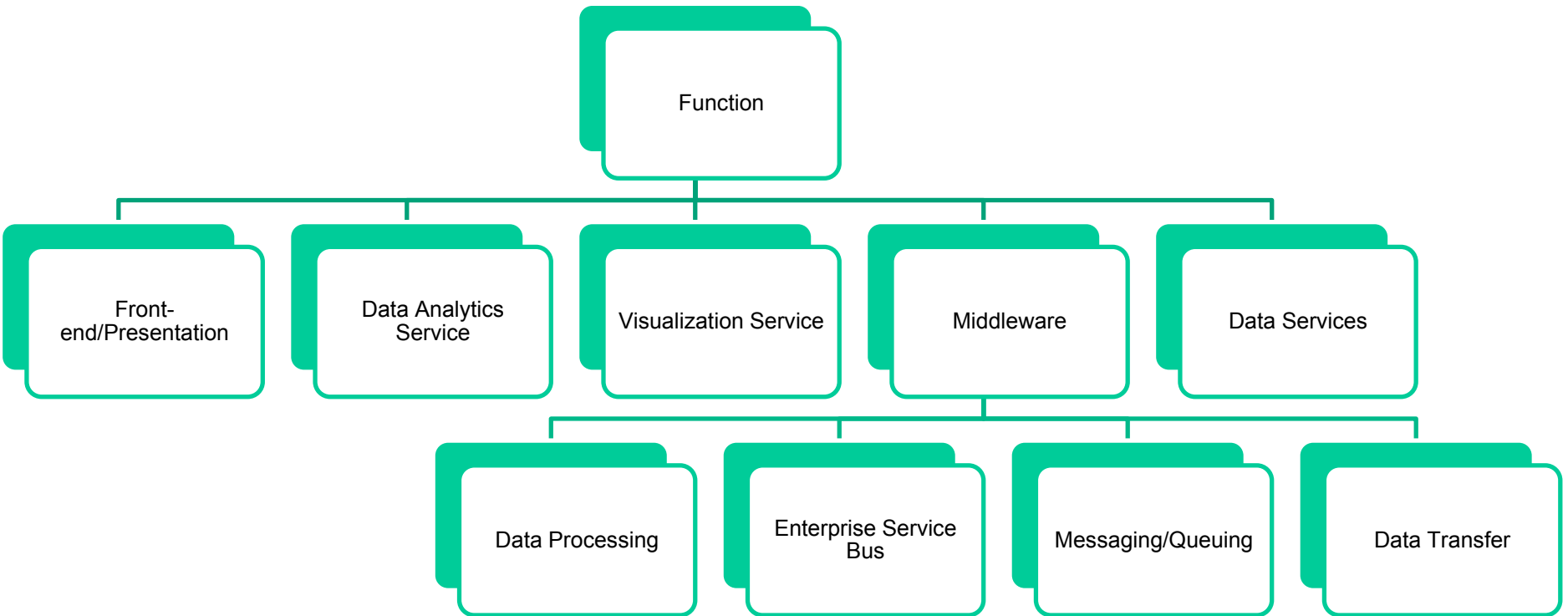> Part = a (composite) service unit

- Which system infrastructures are used?

- Which interfaces  are suitable for units?

- Which programming models are used within units?

- Which are fundamental units to be used?

- How do different units interact?

- Which non-functional parameters are important and how to measure them?

DISTRIBUTED SYSTEMS GROUP

# Fundamental concepts – system infrastructure unit

System infrastructures

- Human-based Computing
- Cloud Services
- Fog/Edge Computing Nodes
- Cluster
- Grid
- High Performance Server
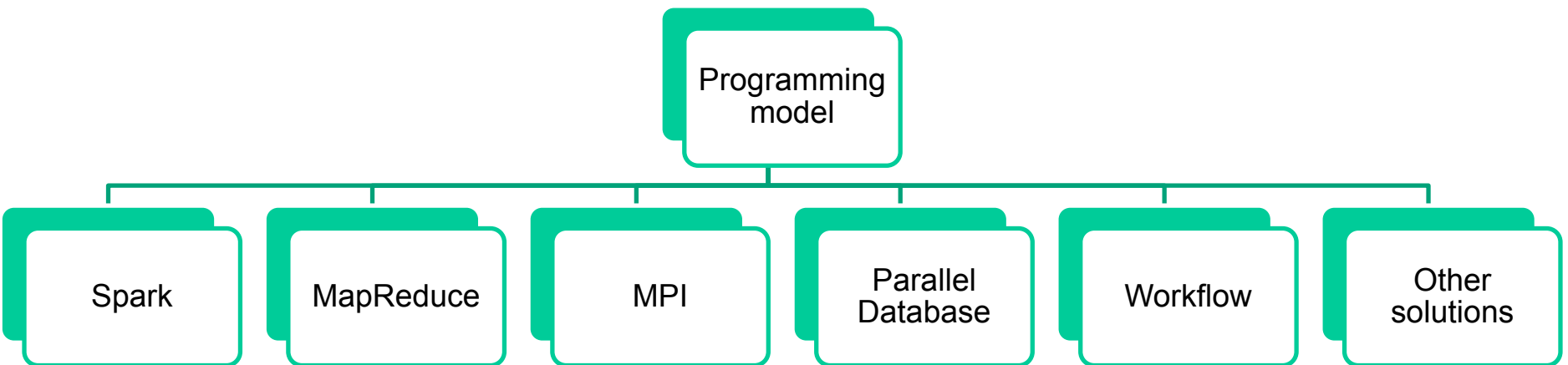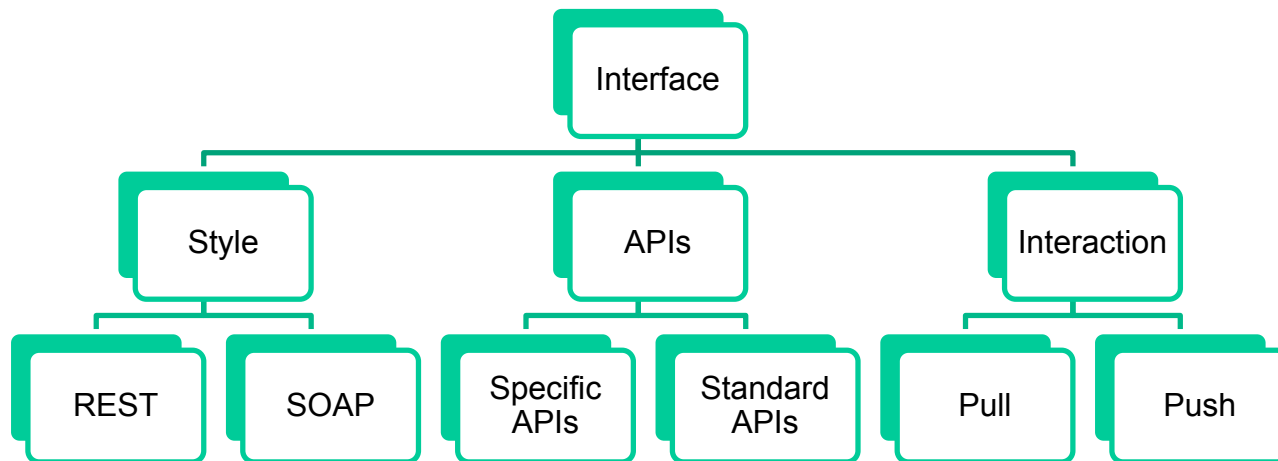
DISTRIBUTED SYSTEMS GROUP
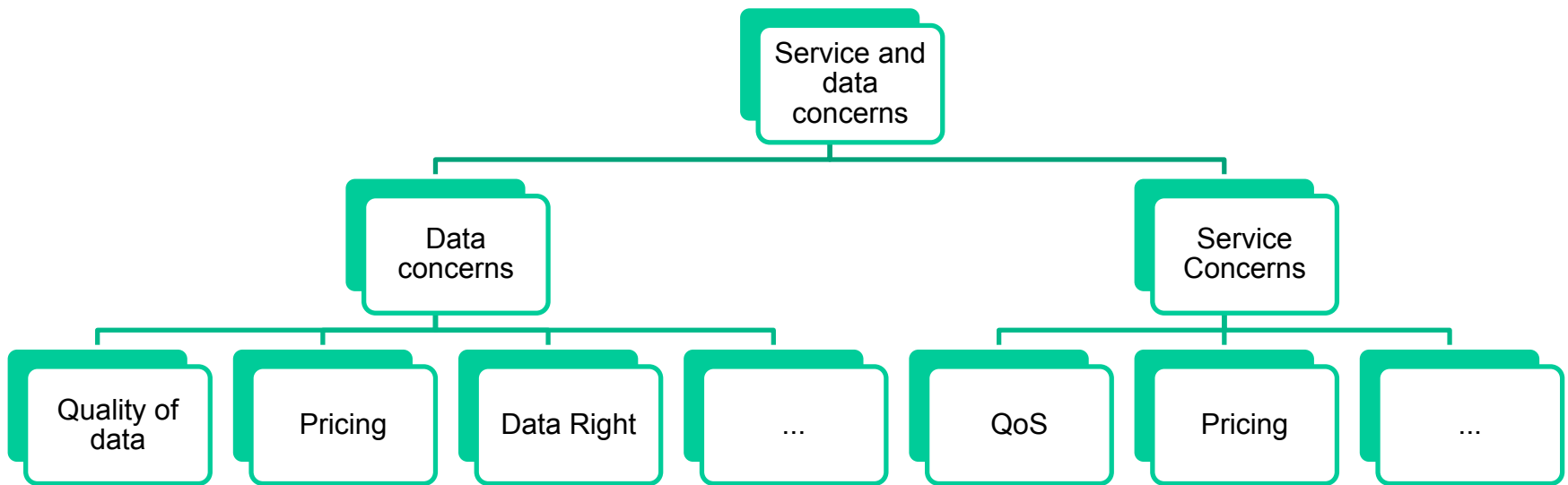
# Fundamental concepts – unit functions

# Fundamental concepts – programming model within units

# Fundamental concepts – interfaces between units
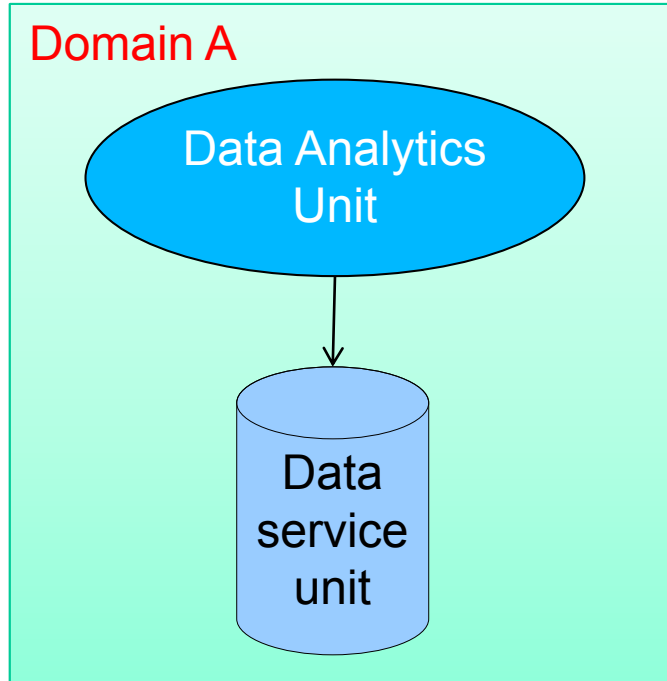
19

# Fundamental concepts – services and data concerns

20

# You see we need to deal with many techniques and frameworks

DISTRIBUTED SYSTEMS GROUP

# WE NEED TO START FROM DATA ANALYTICS WITHIN A SINGLE SYSTEM

DISTRIBUTED SYSTEMS GROUP

# Data analytics within a single system

**Domain A**



- They are complex enough but do not meet all requirements
- In a single domain
  - Tightly coupled computing infrastructures
    - E.g., in the same cloud
  - Computation and data are close
  - Several concerns can be by-passed

Not always provisioned under the „Service Unit" model

DISTRIBUTED SYSTEMS GROUP

# Data analytics within a single system

## Some papers

1. Andrew Pavlo, Erik Paulson, Alexander Rasin, Daniel J. Abadi, David J. DeWitt, Samuel Madden, and Michael Stonebraker. 2009. A comparison of approaches to large-scale data analysis. In Proceedings of the 2009 ACM SIGMOD International Conference on Management of data (SIGMOD '09), Carsten Binnig and Benoit Dageville (Eds.). ACM, New York, NY, USA, 165-178. DOI=10.1145/1559845.1559865
   http://doi.acm.org/10.1145/1559845.1559865

2. Leonardo Neumeyer, Bruce Robbins, Anish Nair, Anand Kesari: S4: Distributed Stream Computing Platform. ICDM Workshops 2010: 170-177

3. Jerry Chou, Mark Howison, Brian Austin, Kesheng Wu, Ji Qiang, E. Wes Bethel, Arie Shoshani, Oliver Rübel, Prabhat, and Rob D. Ryne. 2011. Parallel index and query for large scale data analysis. In Proceedings of 2011 International Conference for High Performance Computing, Networking, Storage and Analysis (SC '11). ACM, New York, NY, USA, , Article 30 , 11 pages. DOI=10.1145/2063384.2063424 http://doi.acm.org/10.1145/2063384.2063424

4. Boduo Li, Edward Mazur, Yanlei Diao, Andrew McGregor, Prashant J. Shenoy: A platform for scalable one-pass analytics using MapReduce. SIGMOD Conference 2011: 985-996

5. Fabrizio Marozzo, Domenico Talia, Paolo Trunfio: A Cloud Framework for Parameter Sweeping Data Mining Applications. CloudCom 2011: 367-374

6. Yingyi Bu, Bill Howe, Magdalena Balazinska, Michael D. Ernst: HaLoop: Efficient Iterative Data Processing on Large Clusters. PVLDB 3(1): 285-296 (2010)

DISTRIBUTED SYSTEMS GROUP

# Data analytics within a single system – some examples

Message Passing Interface (MPI) + Cluster-based File system

Parallel Database (SQL/NonSQL)

MapReduce + Google File System

Yahoo S4

Hadoop + HDFS

Spark

Dryad+LINQ

Scientific/Business Workflow

A short, good overview in Chapter 6: Cloud Programming and Software Environments, Book: Distributed and Cloud Computing – from Parallel Processing to the Internet of Things, Kai Hwang, Geoffrey C. Fox and Jack J Dongarra, Morgan Kaufmann, 2012
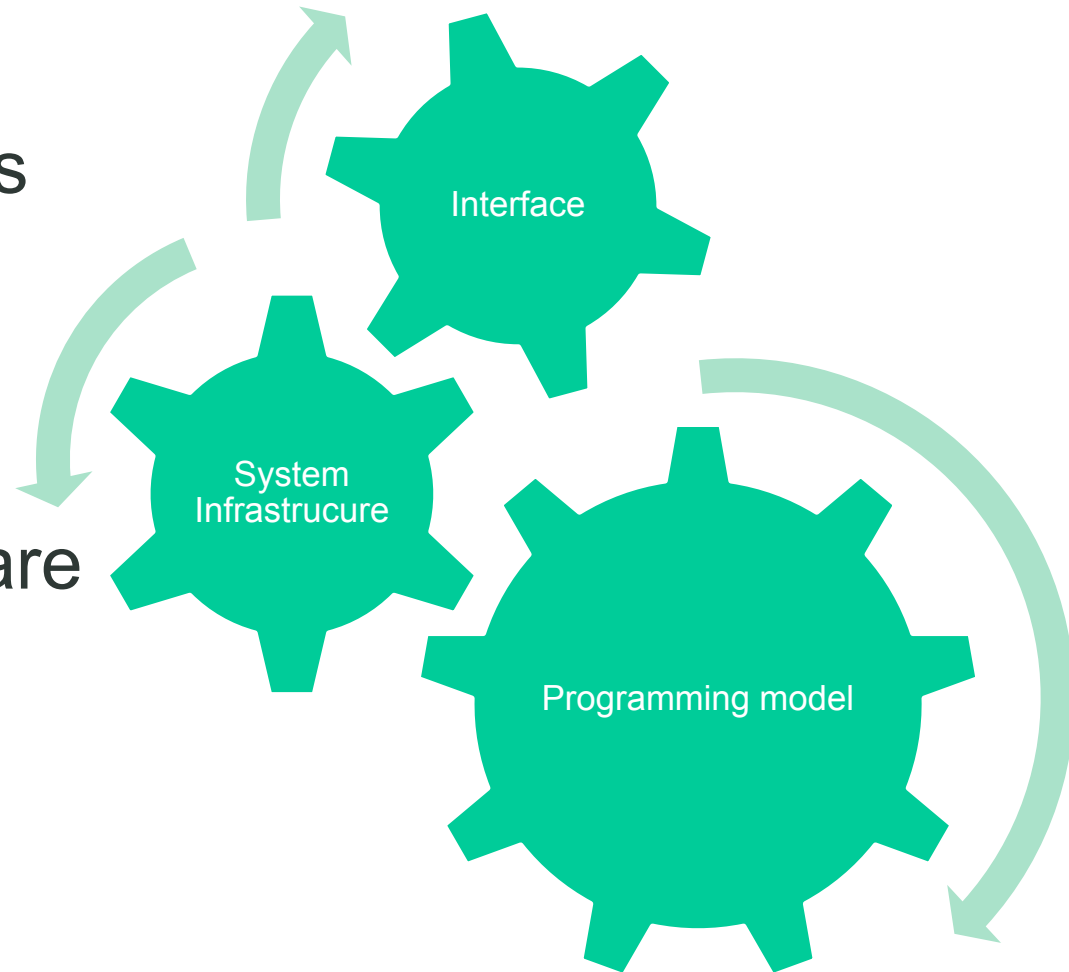
DISTRIBUTED SYSTEMS GROUP

# WHY SHOULD ANALYTICS UNITS BE „CLOSED" TO DATA UNITS?

DISTRIBUTED SYSTEMS GROUP

# WHICH CONCERNS COULD BE IGNORED IN SINGLE SYSTEM DATA ANALYTICS?

DISTRIBUTED SYSTEMS GROUP

**WHICH ARE THE ISSUES THAT WE NEED TO CONSIDER WHEN OUR DATA UNITS ARE IN DIFFERENT SYSTEMS?**
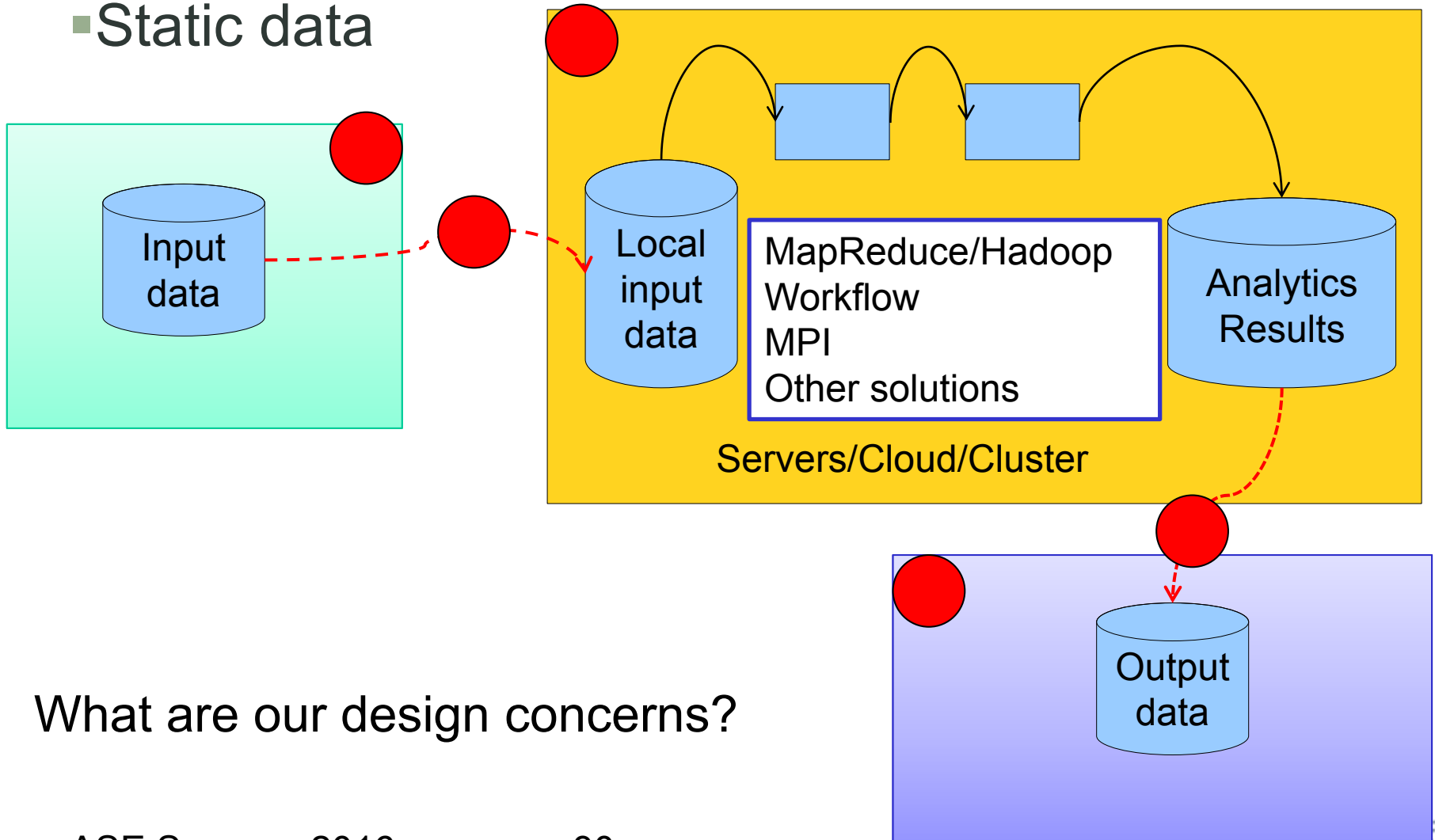
DISTRIBUTED SYSTEMS GROUP

# Data analytics across multiple systems – design choice

- Programming models for data analytics service

- Data service units

- Supporting middleware units

Interface

System Infrastrucure

Programming model

DISTRIBUTED SYSTEMS GROUP

# Data analytics across multiple systems – programming models (1)

■ Static data



```
MapReduce/Hadoop
Workflow
MPI
Other solutions
```

Local input data → Analytics Results

Servers/Cloud/Cluster

Input data → Local input data

Analytics Results → Output data

What are our design concerns?

DISTRIBUTED SYSTEMS GROUP

# Data analytics across multiple systems – programming models (2)

■ Near-realtime data

Stockmarket
Social media
M2M

Input data

Complex event processing
Streaming data analysis
Other solutions

Analytics Results

Servers/Cloud/Cluster

Output data

What are our design concerns?

DISTRIBUTED SYSTEMS GROUP

# Data analytics across multiple systems – programming models (3)

- Near-realtime data



Big data (e.g., satellite images)

Input data

MPI
Workflow
Other solutions

Analytics Results

Servers/Cloud/Cluster

Output data

*What are our design concerns?*

DISTRIBUTED SYSTEMS GROUP

# Data analytics across multiple systems – data service units

**Data Analytics Unit**

Read/write

data

**Cluster file**

## Interface

- Read/write data via direct , low-level read/write via IO

## System

- Cluster or cluster of clusters
- Can be very large

## Programming model
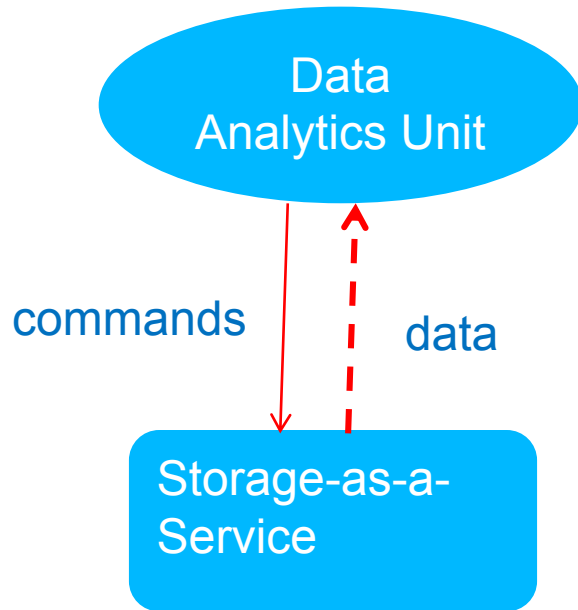
- Usually parallel processing

NFS

Lustre

Hadoop File System

Google file system

DISTRIBUTED SYSTEMS GROUP

# Data analytics across multiple systems – data service units



**Data Analytics Unit**

commands

data

**Storage-as-a-Service**

**Interface**
- Direct data transfer via REST/SOAP APIs
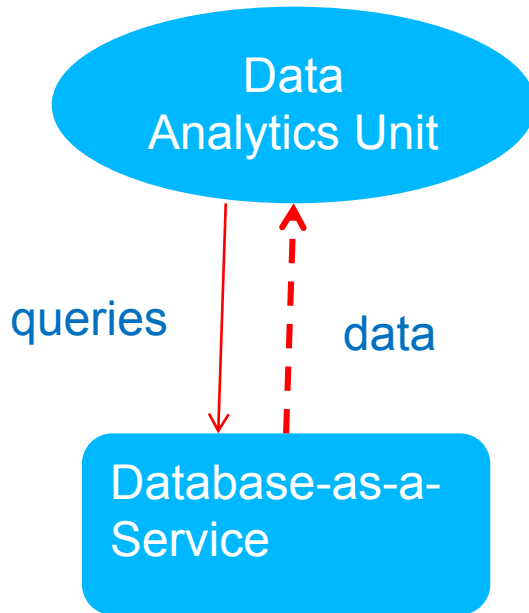
**System**
- Decouple between analytics and storage

**Programming model**
- May require middleware for data transfer
  - Request via SOAP/REST
  - Real data transfer done by external middleware
- A rich set of programming models can be used

Amazon S3
(SOAP/REST API)

Google Storage Service
(REST API)

DISTRIBUTED SYSTEMS GROUP

# Data analytics across multiple systems – data service units

**Data Analytics Unit**

queries

data

**Database-as-a-Service**

**Technology**

## Interface
- REST/SOAP APIs
- Mainly for commands and results

## System
- Decouple between analytics unit and database
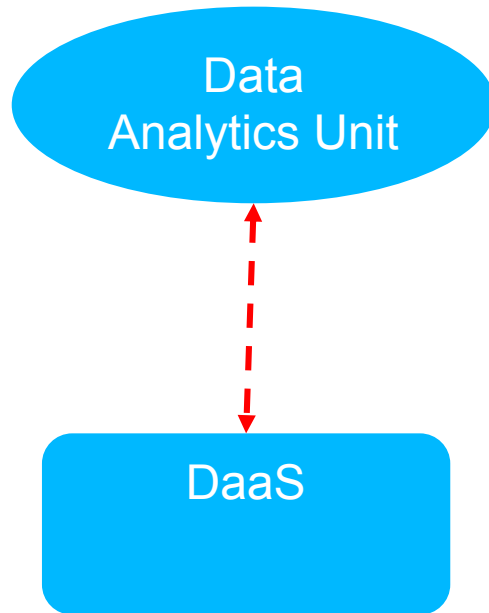- Database as a sevice can be very large

## Programming model
- Analytics can be done at both sides
- Analytic units can use any programming models
- Database-as-a-service can perform  a lot of analytics
  - Parallel database operations

MongoDB/MongoLab
Amazon DynamoDB
Amazon SimpleDB
Cloudant Data

SkySQL
Amazon RDS
Microsoft SQL Azure
Clustrix DBaaS

DISTRIBUTED SYSTEMS GROUP

# Data analytics across multiple systems – data service units

**Data Analytics Unit**

**DaaS**

**Technology**

Infochimps
Microsoft Azure
Xively
GNIP

## Interface

- Data transfer can be uni or bi-direction
- REST/SOAP APIs

## System

- Both systems for DaaS and for analytics units can be very large

## Programming model

- Can be any

DISTRIBUTED SYSTEMS GROUP

# Middleware service unit for transfering large data -- GlobusOnline
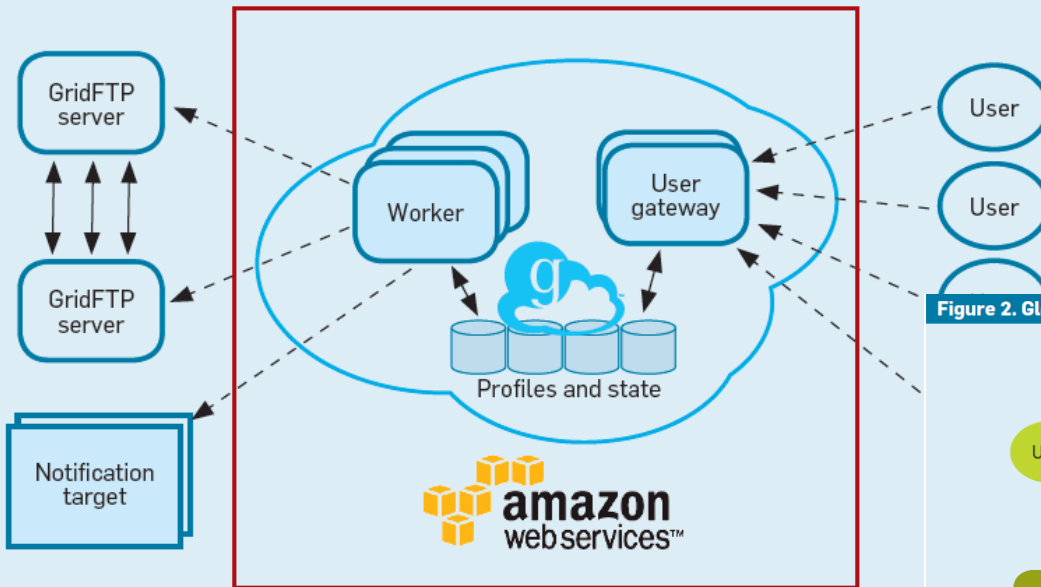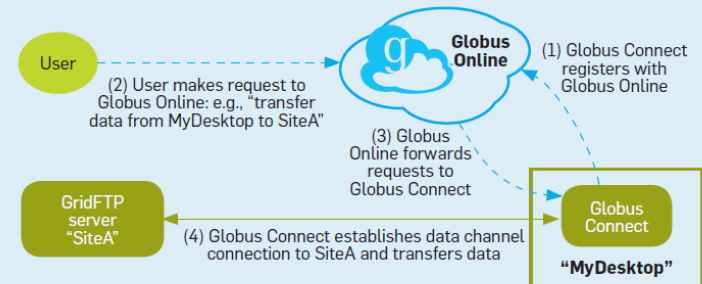


Figure 1. Globus Online architecture.



Figure 2. Globus Connect architecture.

Source: Bryce Allen, John Bresnahan, Lisa Childers, Ian Foster, Gopi Kandaswamy, Raj Kettimuthu, Jack Kordas, Mike Link, Stuart Martin, Karl Pickett, and Steven Tuecke. 2012. Software as a service for data scientists. Commun. ACM 55, 2 (February 2012), 81-88. DOI=10.1145/2076450.2076468 http://doi.acm.org/10.1145/2076450.2076468

ASE Summer 2016

37

DISTRIBUTED SYSTEMS GROUP

# Middleware service units for messages/queuing

- Advanced Message Queuing Protocol (AMQP)
- Simple (or Streaming) Text Orientated Messaging Protocol (STOMP)
- Specific protocols/APIs

| StormMQ | RabbitMQ | Amazon SQS |
|---------|----------|------------|

DISTRIBUTED SYSTEMS GROUP

So many types of services from different providers. Anyway to simplify the management of service units for the developer/operator?
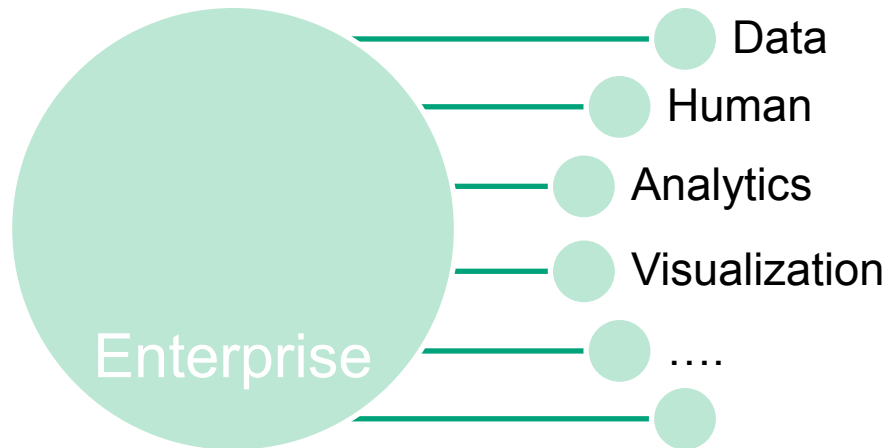
DISTRIBUTED SYSTEMS GROUP

# API MANAGEMENT

# Ecosystem view for advanced service engineering

- Complex data analytics applications → need to understand potential service units from an <span style="color:red">ecosystem perspective</span>

  - Interdependent systems: Social computing, mobile computing, cloud computing, data management, etc.

  - Different types of information are linked

  - Blending vertical and horizontal analytics

  - Different functions (analytics, visualization, communications, etc.)

  - Too many different types of customers (and their interactions)

DISTRIBUTED SYSTEMS GROUP

# APIs

- APIs are key! Why?
  - Enable access to data and function from entities in your ecosystem
  - Virtualization



- An API is an asset
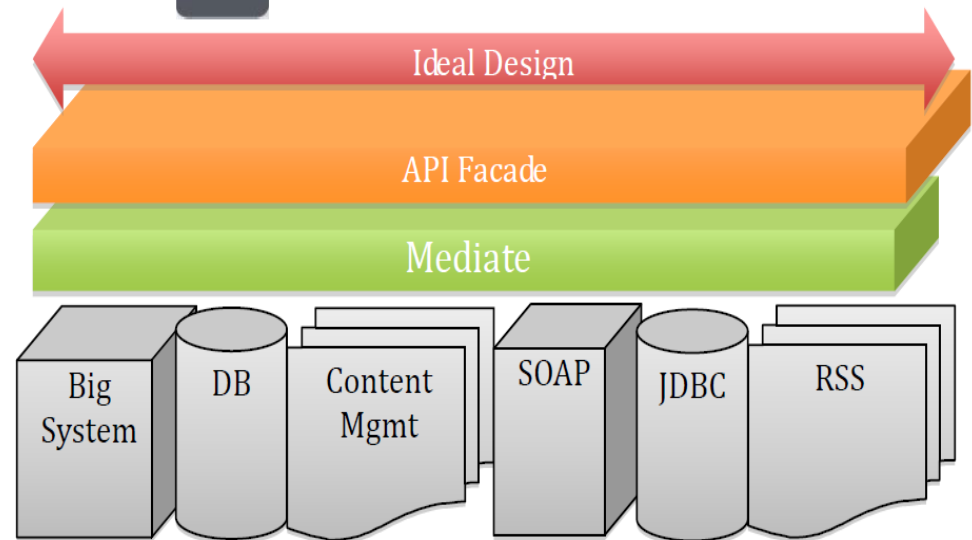  - We need to have lifecycle, pricing, management, etc.

Check http://www.apiacademy.co for some useful tutorials
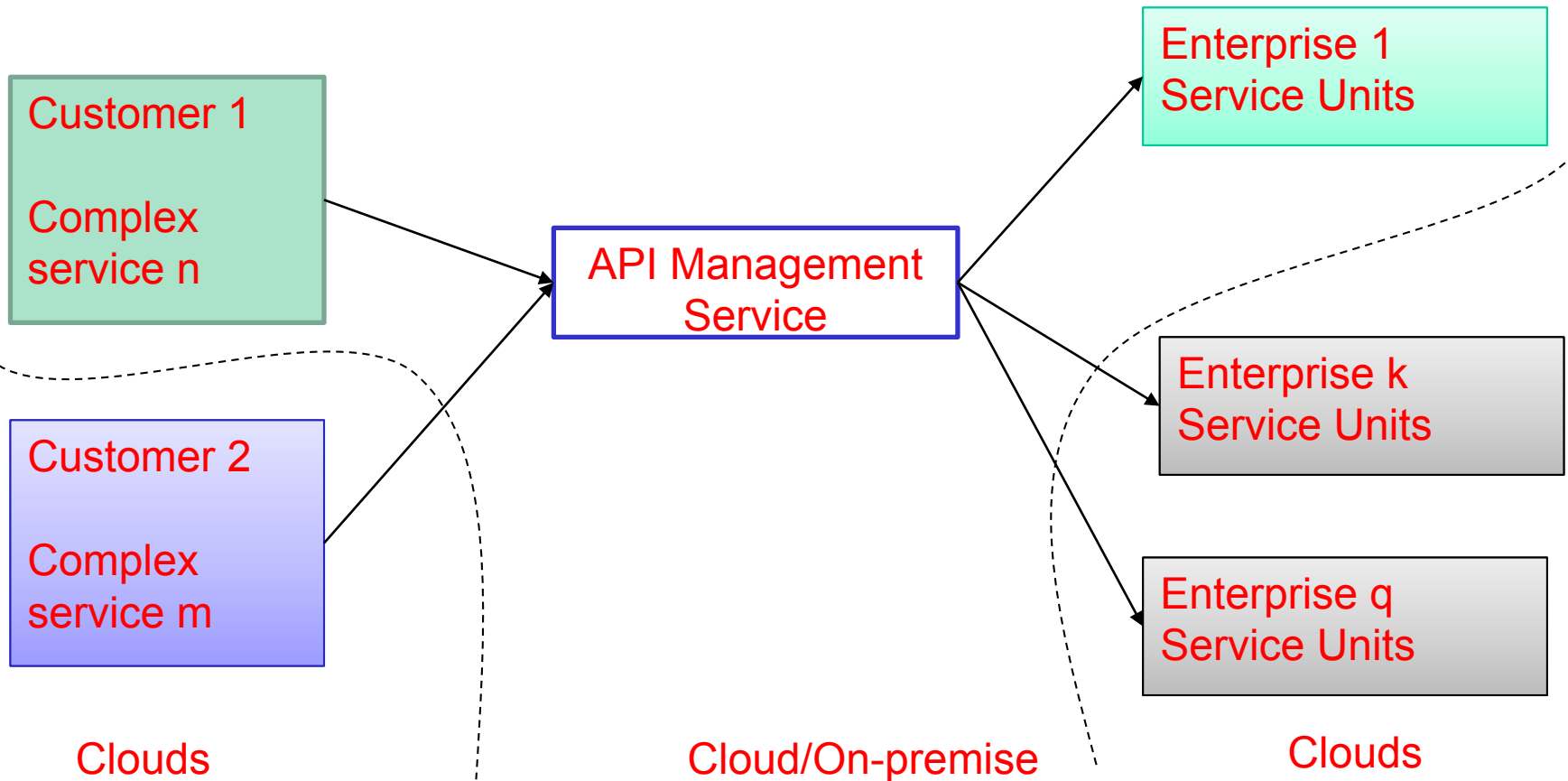
DISTRIBUTED SYSTEMS GROUP

# API Fasade



Sourre:
https://en.wikipedia.org/wiki/Facade_pattern

Source: Web API Design, Brian Mulloy
http://apigee.com/about/resources/ebooks/web-api-design

# API management & APIs as a service

Managing APIs ecosystems

```
Customer 1

Complex
service n
```

```
Customer 2

Complex
service m
```

API Management
Service

Enterprise 1
Service Units

Enterprise k
Service Units

Enterprise q
Service Units

Clouds

Cloud/On-premise

Clouds

DISTRIBUTED SYSTEMS GROUP

# Development of APIs

- Not just the functions behind the APIs
    - This we have learned since a long time
- Emerging (business/service) management aspects
    - Usage control and security
    - Any where from any device for any customer
        - Interfaces (communications, inputs/output formats)
    - APIs as a service:
        - Availability and reliability of APIs are important – think APIs are similar to a service that your client will consume

DISTRIBUTED SYSTEMS GROUP

# Prevent too many accesses?

Client → 100000 requests/s → Service ❌

Client → API Management Service → Service

```
REST_FRAMEWORK = {
    'DEFAULT_THROTTLE_CLASSES': (
        'rest_framework.throttling.AnonRateThrottle',
        'rest_framework.throttling.UserRateThrottle'
    ),
    'DEFAULT_THROTTLE_RATES': {
        'anon': '100/day',
        'user': '1000/day'
    }
}
```

Code: http://www.django-rest-framework.org/api-guide/throttling/#how-throttling-is-determined

ASE Summer 2016

# How can we use API management for data/service contracts?

DISTRIBUTED SYSTEMS GROUP

# Issues on APIs management

- Publish
  - Business and operation planning
    - API usage schemes (e.g., pricing, data concerns)
    - API payload transform policies
    - API throttling
  - API publish and discovery  (like service discovery?)
- Management
  - Management roles in enterprises, versions, etc.
- Monitoring and analytics
  - monitoring and analytics information (availability, types of customers, usage frequencies, etc.)

# Some well-known frameworks

- http://apigee.com

- Oracle API management:
  http://www.oracle.com/us/products/middleware/soa/api-management/overview/index.html

- http://wso2.com/api-management/

- http://www.ca.com/us/lpg/layer-7-redirects.aspx

- https://www.mashape.com/

- http://apiaxle.com/

DISTRIBUTED SYSTEMS GROUP

# Build your own APIs ecosystem

- Which APIs you need?  Which ones are crucial for  you to build complex services?
    - Data APIs
        - Data collection
        - Visualization
        - Analytics APIs
    - Communication
    - Coordination of tasks
- API marketplaces → your APIs
- Using existing API platforms to manage your APIs

DISTRIBUTED SYSTEMS GROUP

# Examples of an API marketplace

51

# Use API Management for your mini project?

# **Exercises**

- Read mentioned papers

- Analyze the relationships between programming models and system infrastructures for data analytics across multiple domains

- Examine http://cloudcomputingpatterns.org and see how it supports data analytics patterns

- Develop some patterns for data analytics across multiple systems

- Setup an API management platform for your work

DISTRIBUTED SYSTEMS GROUP

# Thanks for your attention

Hong-Linh Truong
Distributed Systems Group, TU Wien
truong@dsg.tuwien.ac.at
dsg.tuwien.ac.at/staff/truong
@linhsolar

54