



WALKTHROUGH & STATUS

2013 08 28

S. Yoakum-Stover, Ph.D.
M. Andrew Eick

Institute for Modern Intelligence
& Mission Focus



Move it forward. Make it happen.



CONTENTS

- **WHAT IS THING 5**
- **WHY IT'S CALLED THAT**
- **WHAT PROBLEM DOES IT SOLVE**
- **RELEVANCE TO NGA**
- **GEMS - UNIQUE SOLUTIONS & PERSPECTIVES**
- **CURRENT CAPABILITIES**
- **DEVELOPMENT BACKLOG**
- **RELEASE PLAN AND RETROSPECTIVE**
- **SCHEDULE**

Move it forward. Make it happen.



WHAT IS THING 5?



Move it forward. Make it happen.



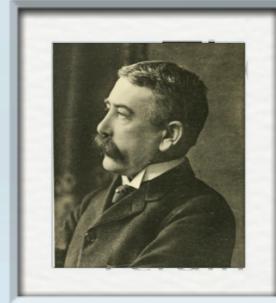
THING 5

A **semiotic** compute & storage apparatus,
built on an open-source **cloud** technology stack,
for the **storage** and **processing** of diverse data at scale,
including geospatial, temporal, human, social, cultural,
behavioral,
as well as traditional Intel data types and
all modalities from documents to streaming video.

Move it forward. Make it happen.

SEMIOTICS

The science which studies how meanings are made and how reality is represented



study of **signs** – anything which stands for something

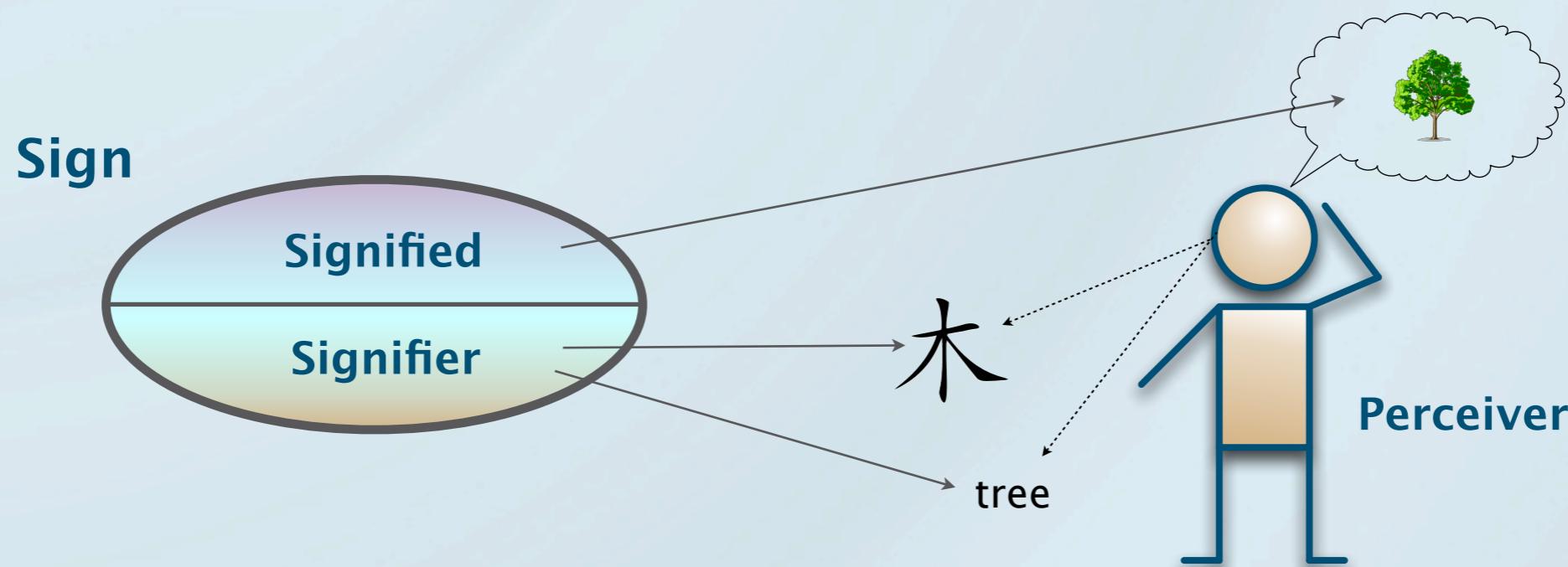
Ferdinand de Saussure (1857–1913)

Developed by

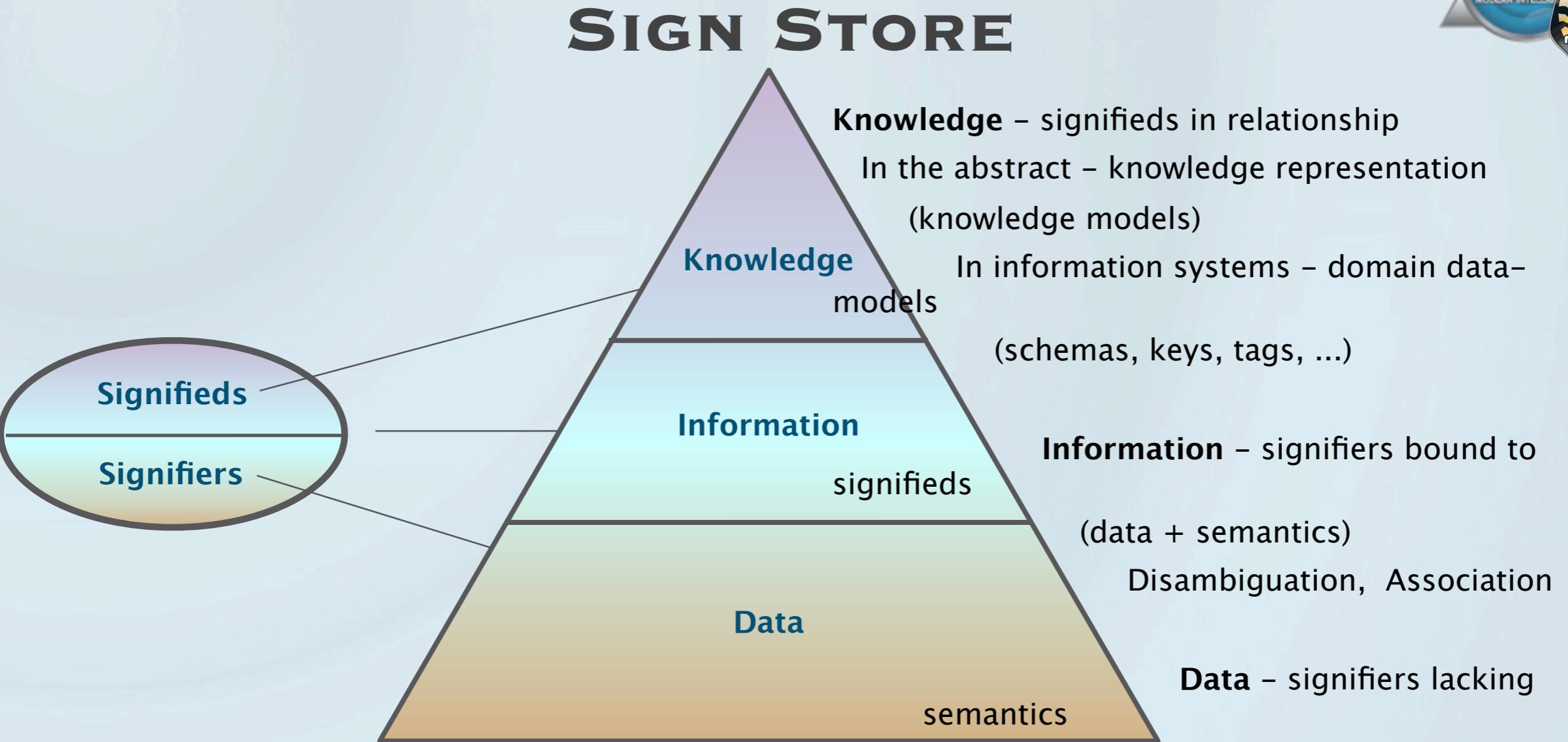
Charles Sanders Peirce (1839)



Logic and Linguistics may be classified as branches of semiotics



To represent a thing, an information system must capture both signifier and signified



In typical information system, signifieds are presented as the storage-model (schema)

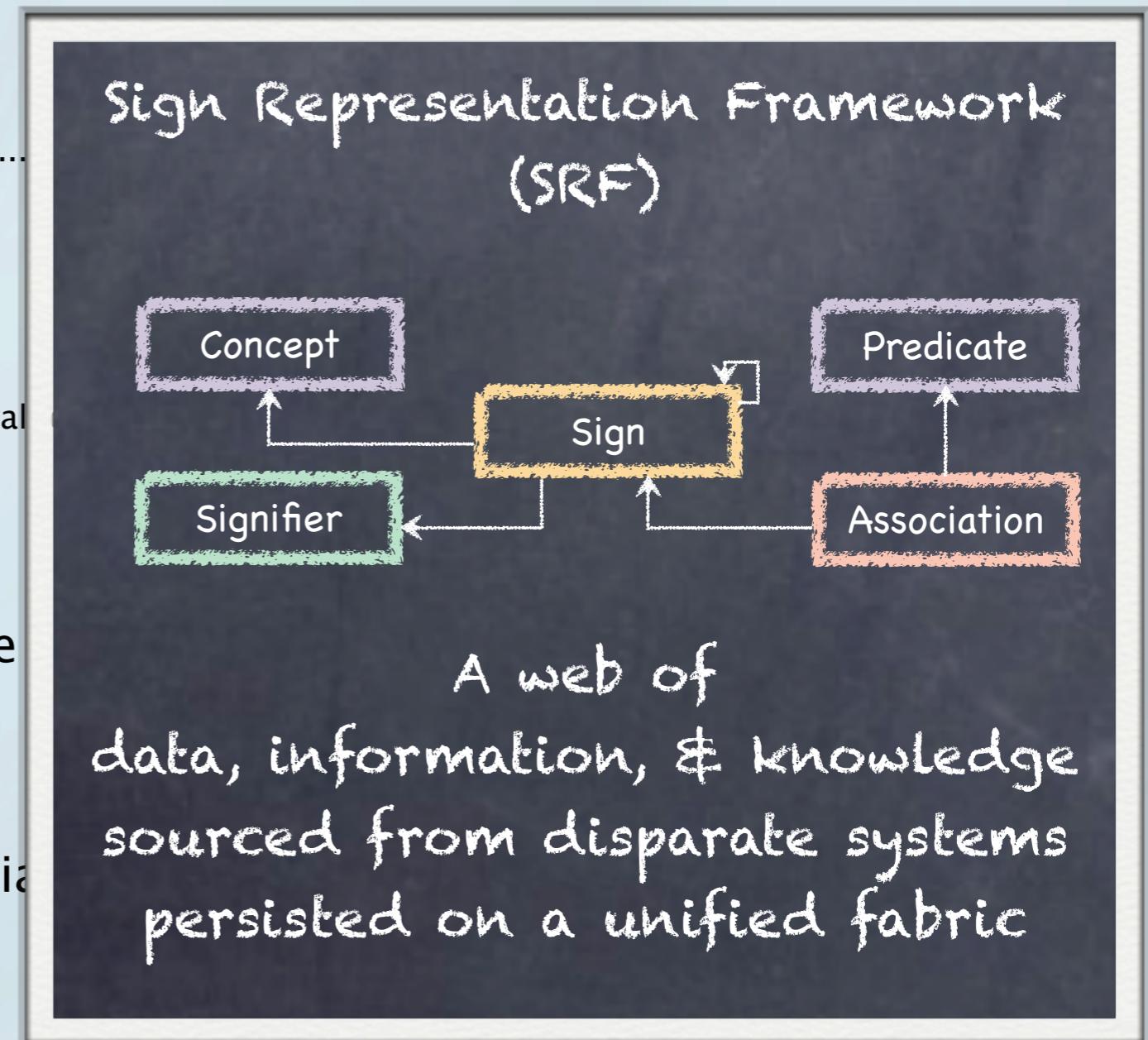
This approach restricts the scope to 1 perspective (set of meanings)

Thing 5 represents signifiers and signifieds as equal citizens
The storage-model is decoupled from the perspective

SIGN REPRESENTATION FRAMEWORK

SIGN STORAGE MODEL FOR ULTRA-LARGE SCALE

- Accommodate all, specify none
 - Any encoding, any modality
 - Video, images, audio, text, numeric series, ...
 - File standards of any kind
 - NITF, MP2TS, XML, IRC, email, HTML, ...
 - Structured data & data-models of all kinds
 - ShapeFiles, spreadsheets, metadata, ...
 - Relational, object, hierarchical, graph, key-val
 - Taxonomies, schemas, ontologies, ...
 - Anything with a geometry / time
 - GIS features, events, ...
- Access all content – unified interface
- Ingest anything
- Preserve domain semantics
- Disambiguate semantically, geospatially, temporally, contextually
- Harmonize data-models, or not
- Integrate information
- Build entirely new kinds of tools, analytics, applications





COMPUTE & STORAGE CLOUD

CHARACTERISTICS

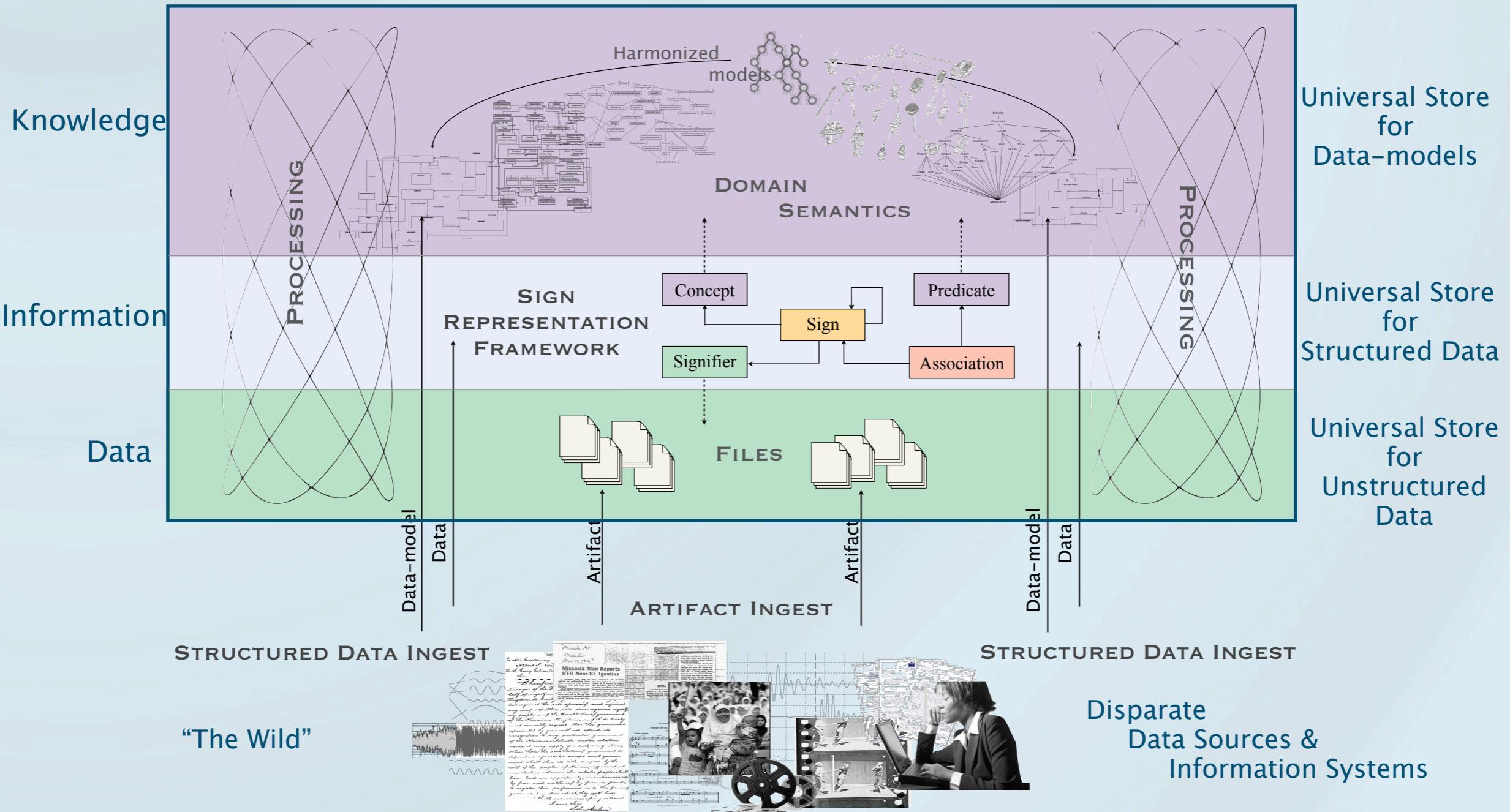
- Internet-scale data store
- Massively parallel computation engine
- Tremendous aggregate bandwidth
- Commodity HW
- Open source SW

ACHIEVING SCALE

- Cloud achieves scale
 - Data size & processing performance
 - Marrying economy with technology
- SRF achieves another level of scale
 - Production and content management
 - Faithfully representing meaning and organizing diversity

SEMIOTIC APPARATUS (SEMAPP)

SIGN-SPACE = SEMIOTIC APPARATUS + DATA + PROCESSING



Move it forward. Make it happen.



BENEFITS & CAPABILITIES

BENEFITS

- No barriers to data ingest, it's fast, simple, and universal
- All users & all processes can access all the data
- Disparate data are unified and operationalized without information loss or distortion
- Richness and meaning of information from the source is preserved
- Data-model harmonization is more powerfully supported, but not required
- True data integration across domains – connecting the dots – is enabled
- Entirely new kinds of analytics, tools, and applications become possible

CAPABILITIES

- Assert, retrieve, delete, tally, query data, information, knowledge
- Search – Keyword, semantic, geospatial, temporal
- Extract, characterize, expose – Make information more discoverable
- Connect – Assert new associations (e.g. social network analysis, registration)
- Surf – Follow associations within and across semantic domains
- Mine – Discover and expose new information (e.g. identification, tracking)

Move it forward. Make it happen.



WHY'S IT CALLED THAT?



Move it forward. Make it happen.



IN THE BEGINNING

CDR KRAFT, SIR... SO WHERE DO WE FIT IN?

I GOT A HANDFUL OF THINGS I WORRY ABOUT

THING 1 IS NVS

THING 2 IS ...

YOU'RE THING 5

YES, SIR :-)



WHAT PROBLEM DOES IT SOLVE?



Move it forward. Make it happen.



THE “DATA” PROBLEM

OUR INTELLIGENCE ASSETS

DATA, INFORMATION, KNOWLEDGE, ANALYTICS, TOOLS, APPLICATIONS...

ARE FRACTURED

IN PHYSICALLY AND SEMANTICALLY DISPARATE SYSTEMS

OPERATIONS ARE IMPEDED

BY OUR INABILITY TO FULLY

SEARCH, EXPLORE, ENRICH, MANAGE, AND EXPLOIT OUR ASSETS

TO CRACK THE DATA / PROCESSING / GLOBAL SENSE-MAKING PROBLEM

THING 5

UNIFIES ALL DATA & SUPPORTS ALL PROCESSING



RELEVANCE TO NGA



Move it forward. Make it happen.



BOLD VISION WICKED BARRIERS

As long as we have the courage to
“**let the data surprise us,**”

I am convinced our analysts will lead to new discoveries—and not only by NGA.
We will **enrich the analysis of others,**
especially the analysis performed by the **all source analyst.**

Letitia Long, Director, NGA



NATIONAL GEOSPATIAL-INTELLIGENCE AGENCY

Know the Earth... Show the Way... Understand the World

NGA Has a *Bold* Vision

Putting the Power of GEOINT in Your Hands



NGA Strategy

Provide online,
on-demand access
to our GEOINT
knowledge

Broaden and deepen
our analytic
expertise to produce
new value

“We will continue to deliver to our varied customer set what they need, when they need it, how they need it. But we have to be thinking about the future. We have to be continually pushing ourselves so that we do remain at the forefront.”

— Letitia A. Long, Director, NGA

Approved for Public Release – NGA Case #13-153

BARRIERS

- The Data Problem
- The Processing Problem
- The Network Effect Problem

THING 5

BREAKS THE BARRIERS



KNOW THE EARTH

MEANS KNOWING WHAT IS - The Physical & the Human terrain

FOUNDATION DATA - From imagery to custom / thematic maps

- Representing our human reality
- Telling the story of the earth and the dynamic things, natural and man-made, that operate upon it
- Foundation for analytic discourse

MAKING & MANAGING FOUNDATION DATA - Core to NGA's core domain

- Ravenous and impatient data consumer
- Starving and languishing amidst troves of data, information, and knowledge
- Something is wrong

THE DATA PROBLEM - Diversity, ambiguity, scale

- How meaning is represented – knowledge representation and semiotics
- Domain models and schemas built into every information system
- Not just more data, profoundly different – really, Really Big Data
- Inability to access, understand, manage and use our data assets as a coherent whole

THING 5

- Accommodates diversity, reduces ambiguity, and handles scale
- Access, understand, manage and use our data assets as a coherent whole
- Feed the production chain with rich, diverse information
- Custom / thematic maps as linked data using information from any source

SHOW THE WAY

MEANS NAVIGATING THAT TERRAIN SAFELY

- ❑ Tactical level of capability that a local understanding enables

ANALYTICS

- ❑ Modeling, simulation, what-if analyses, ...
- ❑ Operationalizing the data
- ❑ Conduct of analytic discourse, a man-machine partnership

MAKING & USING ANALYTICS - Essential to NGA's core domain

- ❑ Scattered across programs and organizations
- ❑ Competing for resources and saddled with overhead
- ❑ Something is wrong

THE PROCESSING PROBLEM - The pipeline model

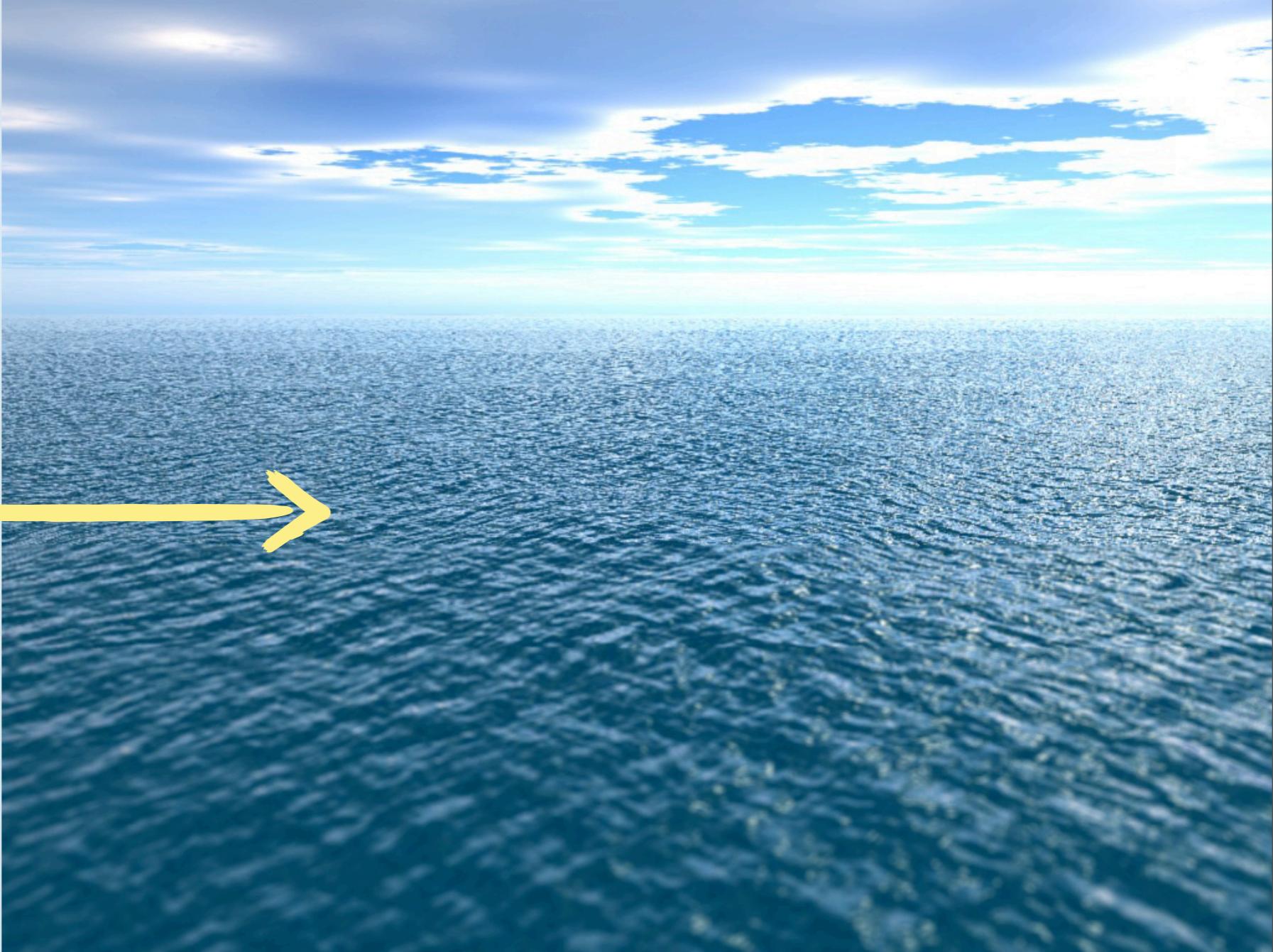
- ❑ Data flowing between processes
- ❑ Data and processing assets are different at every point in the enterprise
- ❑ Bandwidth-limited & complex
- ❑ Inability to leverage, understand, manage and apply our analytic assets as a coherent whole

THING 5

- ❑ Replaces the pipeline model with a wave model – analytic waves over data at rest
- ❑ Leverage, understand, manage and apply our analytic assets as a coherent whole
- ❑ All analytics access and enrich a shared ocean of data, information, knowledge
- ❑ Simple, robust, efficient, powerful, creative



PIPELINES TO WAVES





UNDERSTAND THE WORLD

MEANS KNOWING WHAT IT ALL MEANS

- ❑ A strategic level of capability that a global understanding enables

INTERCONNECTED DATA - Interconnectedness at all scales

- ❑ Putting the pieces together – a living web of information
- ❑ Sharpening the saw – adjusting data and analytic assets in anticipation
- ❑ Product of analytic discourse

MAKING & USING INTERCONNECTED DATA - Essential to NGA's core domain

- ❑ Shared objective across the IC
- ❑ Struggling despite troves of data, information, knowledge, and analytics
- ❑ Something is wrong

THE NETWORK EFFECT PROBLEM

- ❑ Not compounding (each contribution increases the value of the whole and vice versa)
- ❑ No feedback mechanism whereby output (product) enriches the input
- ❑ Inability to fully benefit from achievements and function strategically

THING 5

- ❑ A rich, living web of information naturally emerges
- ❑ Each contribution increases the value of the whole and vice versa
- ❑ Products are an enrichment of the web
- ❑ Fully implement Structured Observation Management
- ❑ Achieve and apply new levels of understanding to achieve mission objectives and make NGA better

PRODUCTS AS INTERCONNECTED DATA

- A product consists of
 - Selected elements of information (signs and associations)
 - Linked together (associations)
 - Displayed by an application to suit a particular purpose
- Examples
 - Notebook
 - Thematic map
 - Social network
 - Geospatial track
- Information elements can appear in any number of different products
- Product creation enriches the web of interconnected data
- Nothing to “re-ingest”
- Every product makes the information more valuable

INSTITUTE FOR MODERN INTELLIGENCE

AN ANALYST NEEDS A GOOD NOTEBOOK

A NOTEBOOK APPLICATION

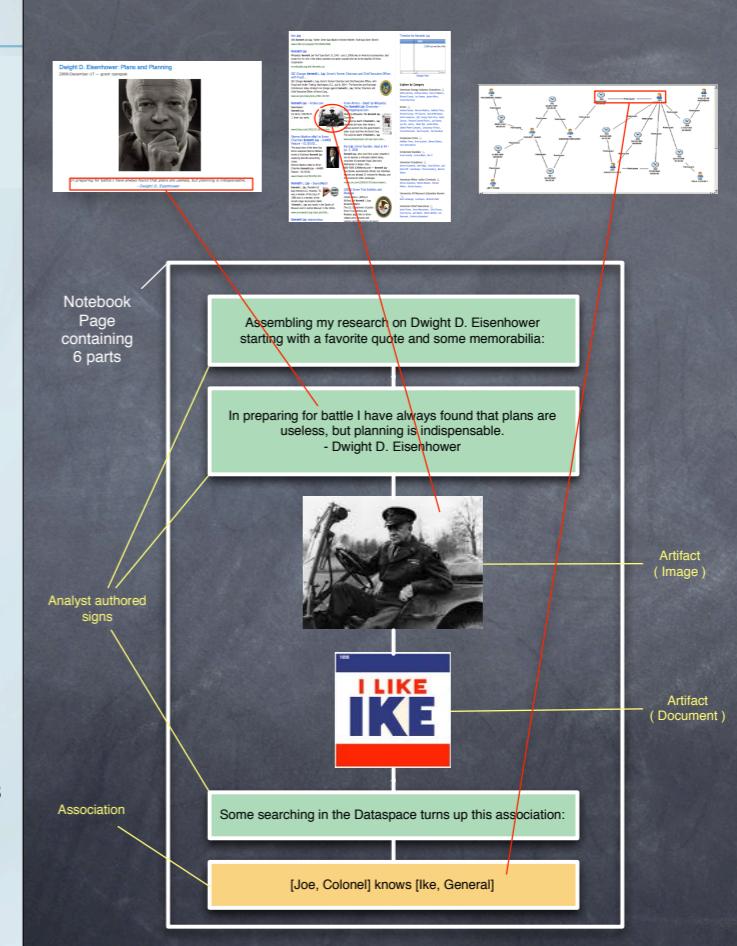
- A tool for performing research and analysis
- Collect, associate, and annotate sign-space elements
 - Artifacts, signs, associations

SIMPLE APPLICATION DATA-MODEL

- A book containing pages
- Each page consists of a linked list of parts
- Any sign-space element can be represented within a part

CONOPS

- Author assembles parts into a page
 - Drag - drop / copy - paste / edit - assert
- This enriches the sign-space with associations
 - No actual document is created
 - Originating author (perceiver) is recorded
- Multiple authors may collaborate on a page



UNCLASSIFIED // PROPRIETARY // COPYRIGHT © 2013 INSTITUTE FOR MODERN INTELLIGENCE, 501(C)3. ALL RIGHTS RESERVED.

10

INFORMATION WEB GROWS AS A SIDE-EFFECT OF ANALYTIC WORK

SPACE & TIME

SPACE & TIME

- ❑ The physics that makes NGA and Thing 5 special
- ❑ The bedrock of the IC
- ❑ The unifying threads throughout our data

THING 5 CORE

- ❑ Semantic disambiguation
 - Binding of a element of **data** (the signifier) to a **concept** (the signified)
 - Forms a Sign – the atom of semantics in Thing 5
 - e.g. [systover@imintel.org, Sender]
- ❑ Geospatial disambiguation
 - A Sign that also has a **geometry** (representing geospatial coordinates) e.g. [[Nile](#), [River](#), [30°10'N 031°06'E](#)]
 - Isomorphic to a GIS feature
 - Applies to associations as well
- ❑ Temporal disambiguation
 - A sign that also has a **period** (representing temporal coordinates) e.g. [[Wall Street bombing](#), [Terrorist Act](#), [19200916](#)]
 - Isomorphic to an event
 - Applies to associations as well

POWER

- ❑ Native GEOINT support
- ❑ Ideal platform for ABI
- ❑ Search, explore, discover, enrich, manage, exploit these unifying threads throughout our data
- ❑ IC solution for GEOINT storage and processing



SEMAPP GEMS



UNIQUE SOLUTIONS & PERSPECTIVES



Move it forward. Make it happen.



UNIQUE SOLUTIONS & PERSPECTIVES

- **A NEW CLASS OF APPLICATIONS (25)**
- **CONTEXT MATTERS (30)**
- **DATA-AS-INFRASTRUCTURE (32)**
- **DATA-MODELS & THE TROUBLE WITH DATA (35)**
- **GEOSPATIAL INDEXING (43)**
- **GEOSPATIAL-TEMPORAL DATA & THE ATOM OF SEMAPP SEMANTICS (49)**
- **INFRASTRUCTURE (51)**
- **INGEST (56)**
- **ORGANIZED TESTING (59)**
- **PARALLEL PROCESS VIDEO (61)**
- **PIPELINES VS. WAVES (63)**
- **TEMPORAL INDEXING (69)**
- **THERE IS NO METADATA (74)**
- **TREATING SEMI-STRUCTURED DATA RIGHT - THE NITF STORY (76)**
- **WHO ARE WE (79)**

Move it forward. Make it happen.



APPLICATIONS



Move it forward. Make it happen.

SEARCH-SURF WEB APP

Artifacts

[Ganesh Set To Act Alongside Amitabh Bachchan](#)
Thaindian.com - May 24, 2010
Ganesh, who made his debut in Tollywood as a Sardarji named Joginder Singh in Abhiyum Naanum, will now play the role of Big B's son in his next film **Kandhar** ...
[Ganesh Set To Act Alongside Amitabh Bachchan](#) - Real Bollywood (blog)
[all 7 news articles »](#)

[NATO installs satellite system on Pak-Afghan border](#)
SAMAA TV - 5 days ago
... NATO has installed the satellite system on Pak-Afghan border prior to launching a massive operation against Taliban in Kandahar province of Afghanistan. ...
[Nato installs surveillance system near Chaman border](#) - The News International
[NATO installs surveillance tower near Chaman border](#) - Daily Times
[all 7 news articles »](#)

Concepts

Category	Count
Person	125
Place	100
State	92
Vehicle	56
Selector	34
Country	5

Sources

Source Type	Count
M3	140
Enron	112
LSIE	110
Imagery	97
Signal	89
Video	55

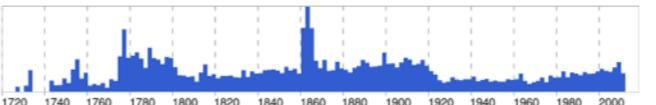
Imagery





A busy street in Photo: Kandahar Gate likely in Kandahar

Timeline

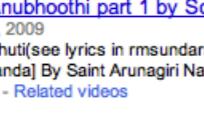


1720-2010 Search other dates

Motion Imagery



[Kandhar People gathering for vote rigging ...](#)
11 min - Sep 8, 2009
[youtube.com - Related videos](#)

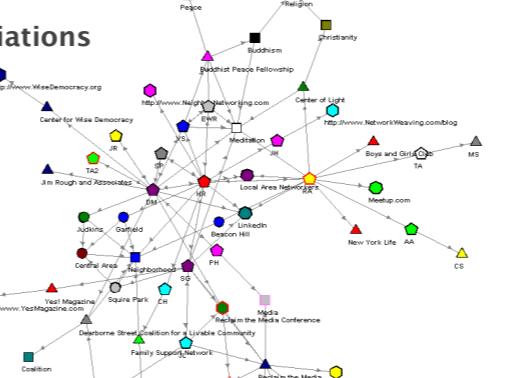


[Kandhar Anuboothi part 1 by Soolamangalam sisters](#)
7 min - Apr 2, 2009
Kandar Anubuthi(see lyrics in rmsundaram1948.blogspot.com [Spiritual Union with Kanda] By Saint Arunagiri Nathar [God Kanda (or Skanda ...
[youtube.com - Related videos](#)



[Taliban Getting Strong in Their Stronghold- PakNewz.com](#)
2 min 32 sec - Nov 27, 2008
getting strong day by day in Kandhar.. Watch Video about
Taliban,NATO,Kandhar by Metacafe.com ... Taliban Taliban NATO
Kandhar ...
[metacafe.com - Related videos](#)

Associations



A network graph illustrating connections between entities such as Center for Wise Democracy, Buddhist Peace Fellowship, Center of Light, Boys and Girls Club, Meetup.com, AA, CS, and many others.

Notes & Actions

I'm working to answer the question: What position connection might Kandahar have to the AAF?

All expresses the following claim:

Load the global news for information and intelligence analysis. This is a great way to keep up with the latest news, politics, and governance developments across the world.

This is the best known photo of Kam Lye.

These resources indicate some interesting connections:

Note cards are showing quick, but there's something here worth checking.

Global news

Opinion

Opinion

Opinion

Opinion

Opinion

Opinion

Opinion

Opinion

Opinion

Load the global news for information and intelligence analysis. This is a great way to keep up with the latest news, politics, and governance developments across the world.

This is the best known photo of Kam Lye.

These resources indicate some interesting connections:

Note cards are showing quick, but there's something here worth checking.

Maps



A SIGN-SPACE NEEDS A GOOD SEARCH ENGINE

- Familiar Google-like interface
- Ask questions of the data
- Boolean and semantic search over the entire Sign-Space
- Navigate hyperlinks to surf associations
- Drill down to originating artifacts and sources
- Quickly locate rich analysis products
 - Notebook pages
 - HorseBlankets
 - Enhanced audio, imagery, and video

Mouse-over for

Move it forward. Make it happen.

UNCLASSIFIED // PROPRIETARY // COPYRIGHT © 2013 INSTITUTE FOR MODERN INTELLIGENCE, 501(C)3. ALL RIGHTS RESERVED

26

Thursday, August 29, 13

AN ANALYST NEEDS A GOOD NOTEBOOK

A NOTEBOOK APPLICATION

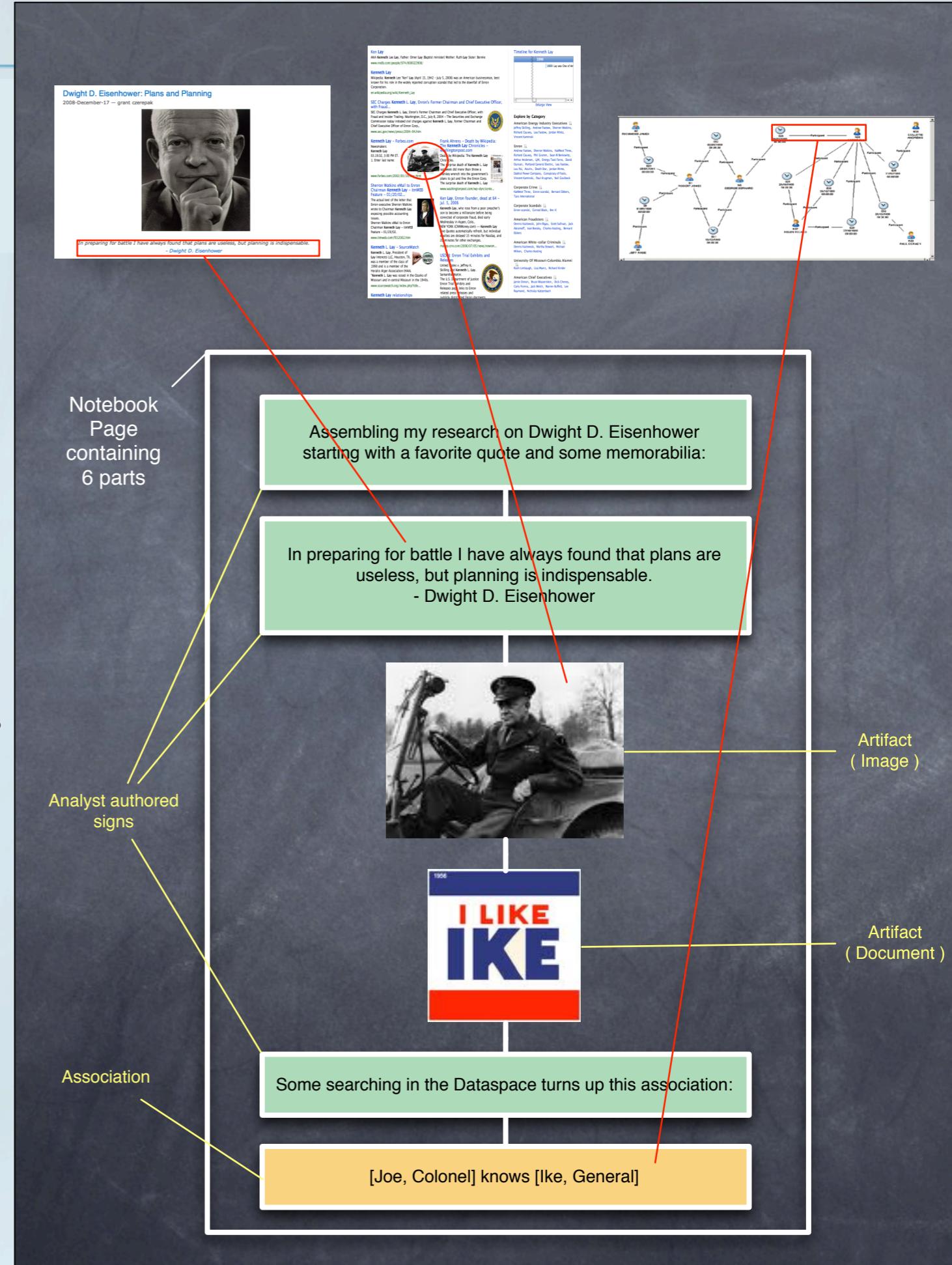
- A tool for performing research and analysis
- Collect, associate, and annotate sign-space elements
 - Artifacts, signs, associations

SIMPLE APPLICATION DATA-MODEL

- A book containing pages
- Each page consists of a linked list of parts
- Any sign-space element can be represented within a part

CONOPS

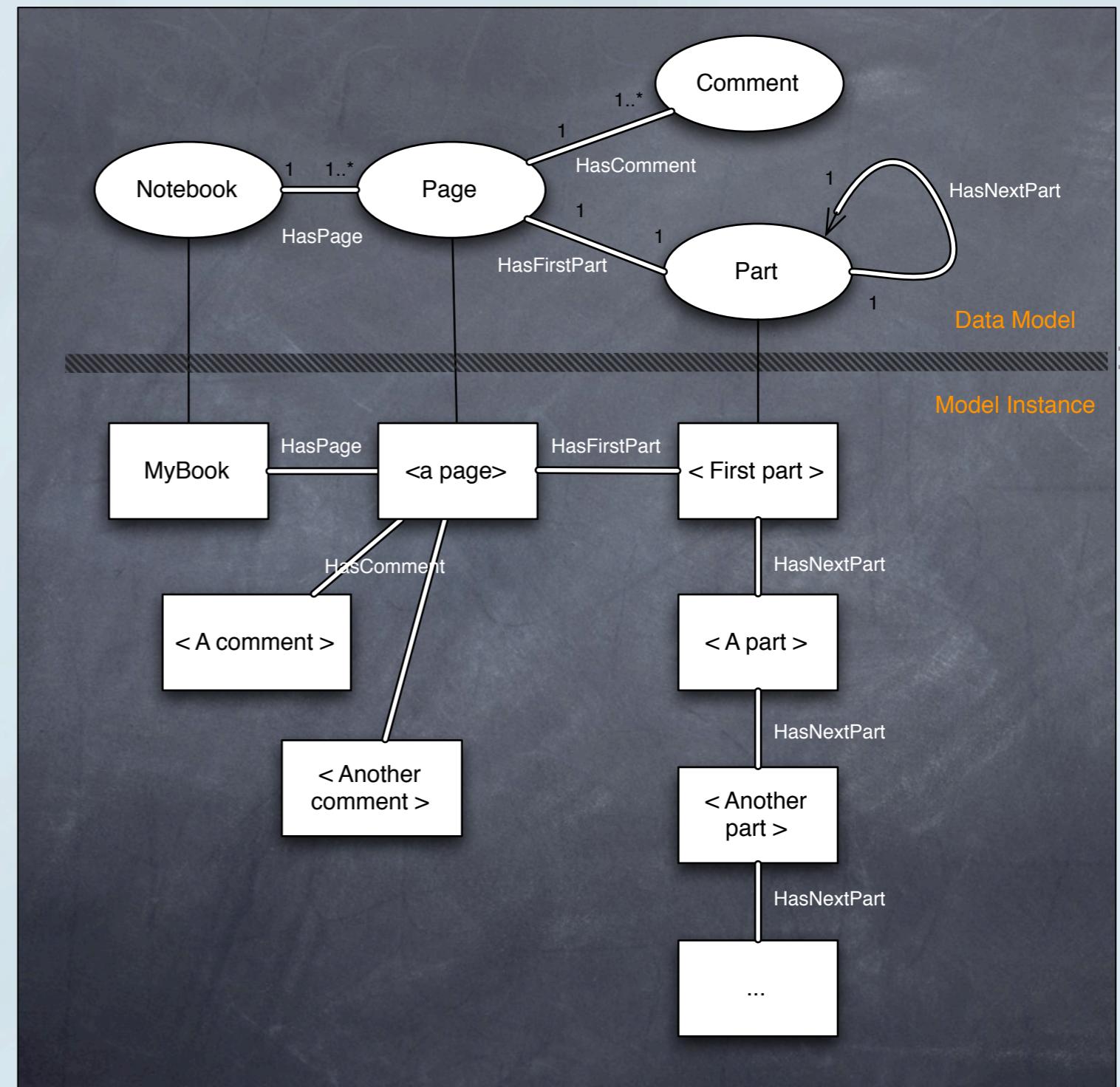
- Author assembles parts into a page
 - Drag – drop / copy – paste / edit – assert
- This enriches the sign-space with associations
 - No actual document is created
 - Originating author (perceiver) is recorded
- Multiple authors may collaborate on a page



NEW SPACE TO WORK IN

APPLICATION DATA-MODEL

- A “Container” Model
 - Imposes minimal semantics
 - Parts serve as containers for sign-space elements
 - Elements in containers carry their original domain-semantics
- A Mash-Up of Semantic Elements
 - Not data-model harmonization
 - Not exactly data integration
 - More than a visual mash-up
 - A persistent, information mash-up
- New Species of Application
 - Enabled by the unified interface
 - Operating on information absent any data-model harmonization
 - A new space in which to work with data!



SIGN-SPACE EVOLVED VIDEO APPLICATION

ENABLES AN ANALYST TO ENRICH VIDEO WITH

- Other contextual data from the Sign-Space
- Annotations including text, images, voice, and telestration

ARCHITECTURE VERY SIMILAR TO SIGN-SPACE NOTEBOOK

- A container model with parts associated by time
 - SSEVA Stream has a MainStream and PartStreams
 - PartStreams have start, length, and placement relative to MainStream
 - PartStreams may contain any Sign-Space element
- Dropping parts into a SSEVA flow automatically creates PartStreams relative to the MainStream
- The SSEVA engine knows how to displays SSEVA streams





CONTEXT MATTERS



Move it forward. Make it happen.

CONTEXT

THE INFORMATION ENVIRONMENT IN WHICH A THING MAKES SENSE

- Context is a kind of Thirdness (Aristotle, Kant, Hegel, Pierce, Husserl, Whitehead, ...)
 - Triad of primary categories of the kinds of things that are
 - Firstness – Existence independent of anything else
 - Characterized by intrinsic qualities (e.g. systover, Mission Focus, D&B)
 - Secondness – Being relative to something else
 - Characterized by dyadic, directed relations between entities (e.g. mother, attorney, wife, employee)
 - Thirdness – A mediation whereby a multiplicity of entities are brought into relation, creating new relationships
 - Characterized by a nexus of related entities (e.g. motherhood, legal system, marriage, business enterprise)
- Examples: Organizational, cultural, economic, religious, technical, ...

CONTEXT IN INFORMATION SYSTEMS

- A data-model implicitly defines a domain context
- Beyond the data-model, context is almost never represented
- Using information outside it's originating domain without preserving the context is perilous

CONTEXT IN THE SIGN-SPACE

- SRF captures firstness, secondness, and thirdness
- Context is represented at several levels
 - Structured data retain their native domain-model associations
 - Data retain their source context (source metadata)
 - Information extracted from artifacts retains the artifact context (mentions)
 - Context may be explicitly represented as signs in the SRF
 - Contextual models may be explicitly represented in SRF



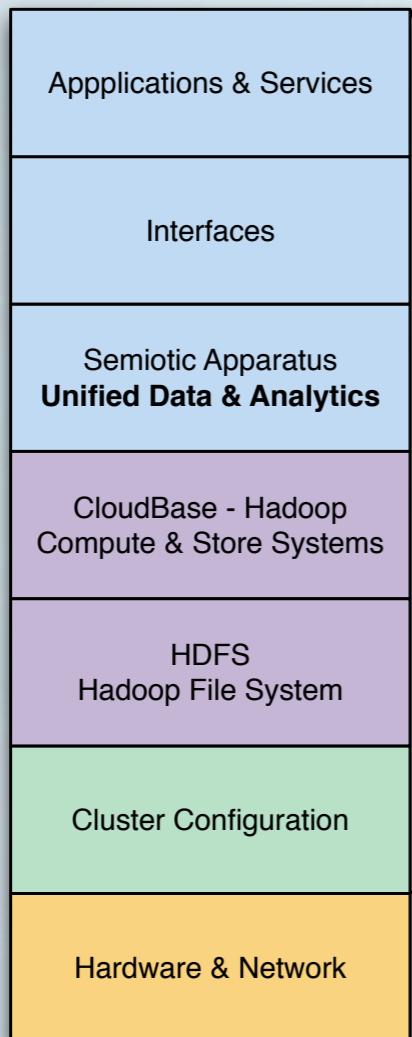
DATA-AS-INFRASTRUCTURE



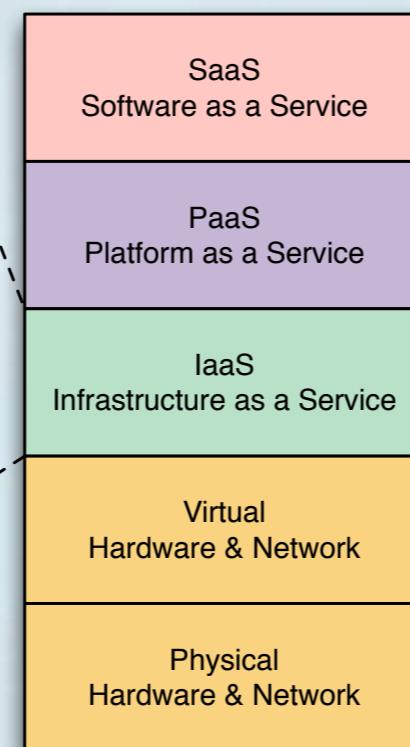
Move it forward. Make it happen.

NEW CAPABILITY → NEW CAPACITY

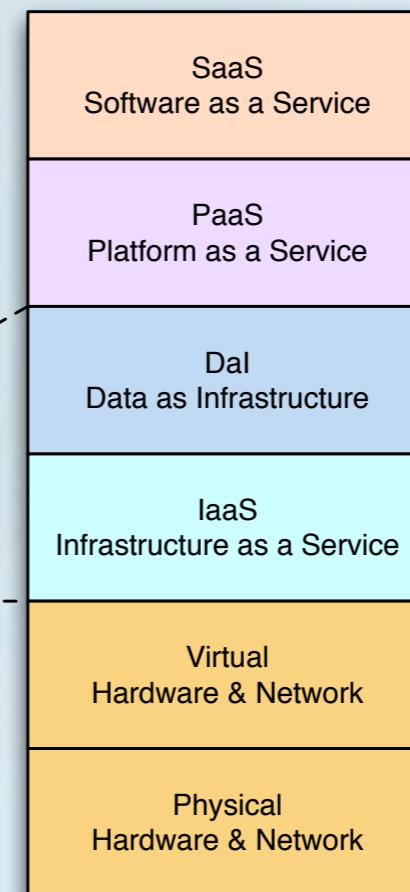
NEA
Compute & Storage
Cloud



Utility Computing
Cloud



Global
Compute Enterprise*
Cloud



Incorporate →
into the utility computing stack

→ Yields
infrastructure / **data** / platform / software capacity



NEW CAPACITY FROM NEW CAPABILITY

**BROADEN NGA'S VISION OF UTILITY COMPUTING BY PROVIDING
"DATA-AS-INFRASTRUCTURE"
(DATA & PROCESSING AS A UBIQUITOUS RESOURCE)**

□ How

- Incorporate the SemApp into the infrastructure layer of the GCE
 - Hook platforms to the Thing 5 cloud at the infrastructure layer
 - Expose cloud data and processing via APIs that hide the underlying transport and storage mechanisms

□ Result

- SemApp becomes an open commodity accessible by default
- Using and contributing to the unified data assets becomes easy by default

**As part of the Data-as-Infrastructure strategy,
Thing 5 transforms the "data problem" into a National Asset
&**

**Achieves the IC's goals for data sharing
as a matter of course**



DATA-MODELS & THE TROUBLE WITH DATA



Move it forward. Make it happen.

DATA-MODELS

SPECIFY

- Vocabulary, Structure, Semantics, and Constraints

BUILT FOR A PURPOSE

- Make a particular analytic, tool, application work efficiently
- Entail a particular perspective on the data
- They're all different for a reason

PROVIDE THE INTERFACE TO AND PERSPECTIVE ON THE DATA

- All interactions with data are mediated by a data-model
 - We go thru the data-model to ask questions about the data
 - What is the **name** of the **Student** with **sid** 53666?
 - How many **Students** have **age** = 18?
 - What is average **gpa** of **Students**?
 - We can not directly question the data
 - What is "Jones"?
 - Is there anything on "Guldu" ?
 - What is "53831" related to?
 - Indexes are no exception
 - Every index reflects a hidden data-m
 - The "things" being indexed
 - e.g. Index on **Students name**

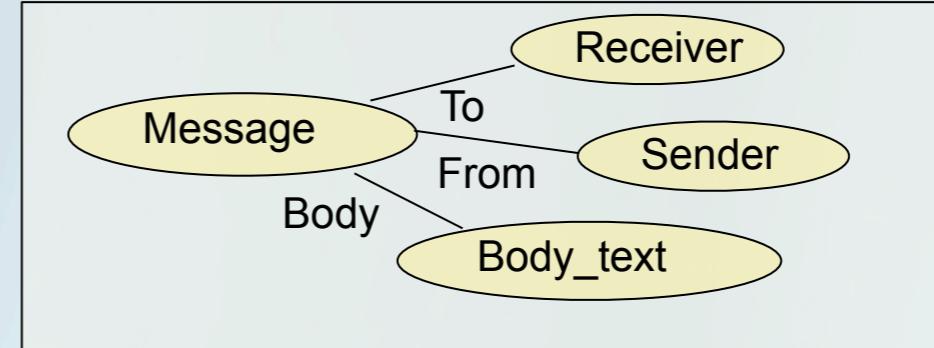
Students

<i>sid</i>	<i>name</i>	<i>login</i>	<i>age</i>	<i>gpa</i>
53666	Jones	jones@cs	18	3.4
53688	Smith	smith@ee	18	3.2
53650	Smith	smith@math	19	3.8
53831	Madayan	madayan@music	11	1.8
53832	Guldu	guldu@music	12	2.0

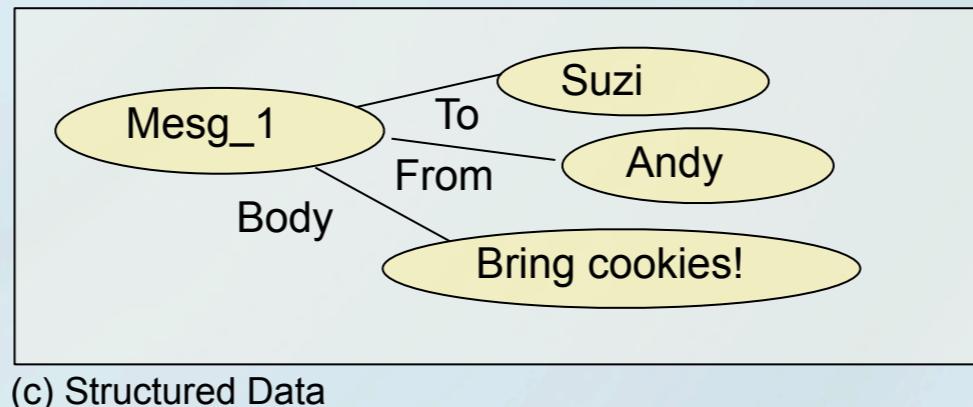
DATA-MODEL IN ACTION (RDB-STYLE)

07/04/07
Suzi,
Bring cookies!
Andy

(a) Unstructured Data



(b) Data-model



(c) Structured Data

The data-model
is imposed
on the data store
(schema)

The data values
are frozen
beneath it (e.g.
We can't say new
things about
Andy)

Message	Receiver	Sender	Body_text
Msg_1	Suzi	Andy	Bring cookies!
...			

(d) Typical database structure

Introducing data
from other sources
is hard as data-
models must
conform to the
storage structure

Once you do that
however, the data
is “integrated” (but
according to only 1
perspective)

Once loaded, it's
painful to modify.

Information is
typically lost /



THE PROBLEM WITH DATA

is...

It Has to Go Somewhere and often goes Lots of Places

Then it Tends to

Stick where it Lands and take the Shape of its Container

In other words,

The Data-Model gets

Imposed on the Data Store and the Data is then Frozen Into it

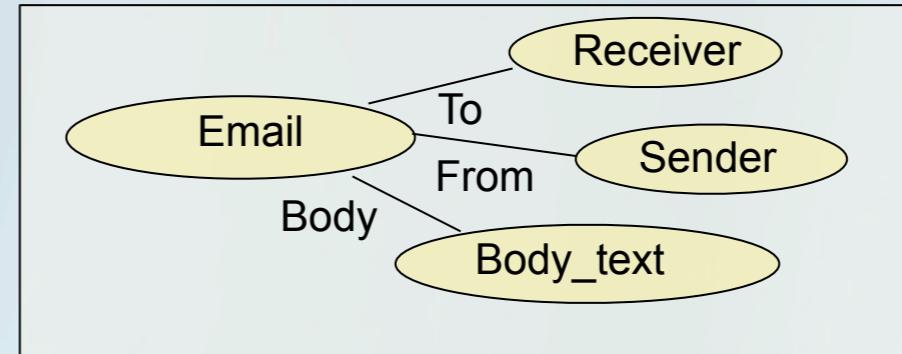
So Don't do That

Move it forward. Make it happen.

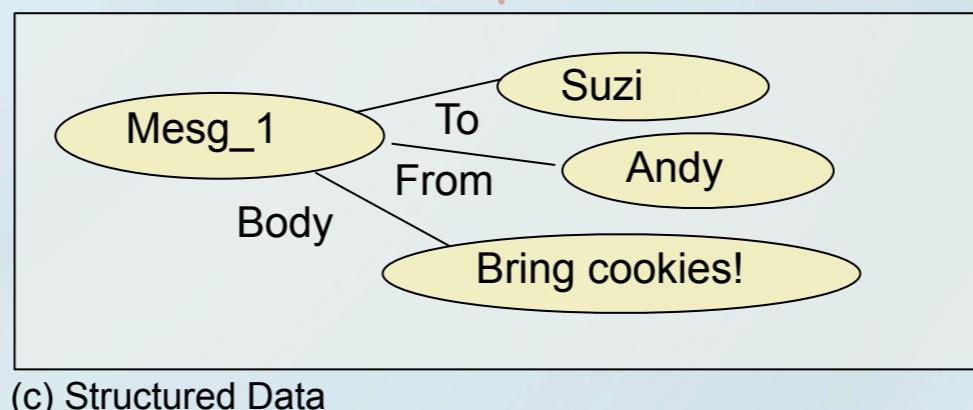
DATA-MODEL IN ACTION (CLOUD-STYLE)

07/04/07
Suzi,
Bring cookies!
Andy

(a) Unstructured Data



(b) Data-model



(c) Structured Data

The data-model is still there, but it's mingled with the data

We still need to know what it is to access our data

We also need to know how data

Introducing data from other sources is easy

Keeping track of their data-models (interfaces) is hard

The data is not integrated

Row	Family	Qualifier	TimeStamp	Value
Mesg_1	Email	Sender	20090224123	Andy
Mesg_1	Email	Body_text	20090224125	Suzi, Bring cookies!

(d) Typical CloudDB structure (e.g. BigTable)



THE PROBLEM WITH CLOUD

is...

The Data Structure has to **GO SOMEWHERE**
AND CAN GO

almost Anywhere
So

Data flows Easily **INTO** its Container
BUT LOSES ITS SHAPE

In other words,

The Data-Model gets
Mixed with the Data in a Specific but **Arbitrary** way

So **Don't** do That

DATA INTEGRATION CONUNDRUM

SITUATION WE HAVE:

A data-model Provides

The Interface to the data AND the Semantics of the data

For data to be integrated it must have a Unified interface

Unified Interface ⇒

A single data-model

A single meaning / perspective

But there is no single right way to represent all Knowledge ⇒

There can be no single data-model for all intel

SITUATION WE WANT:

**A unified interface to All the data
that**



THE SOLUTION

CONSIDER DATA-MODELS FROM A

Higher Level of Abstraction

Distill from there a

Minimal set of Elements Sufficient to Capture
Any data-Model

and then

Build Storage Model based On that That

Move it forward. Make it happen.



GEOSPATIAL INDEXING



Move it forward. Make it happen.



GEOHASH OVERVIEW

**BINARY DECOMPOSITION SCHEME REPRESENTING LAT-LONGS BY RECURSIVELY DIVIDING THE ASSOCIATED ANGULAR RANGES
([-90, 90], [0, 180])
IN HALF**

Bit value of 0 indicates the lower/left half. Bit value of 1 indicates the upper/right half

01111100000000 (long), **101111001001** (lat)

The two lat-long bit sequences are interleaved from left to right with even bits taken for the longitude

01101 11111 11000 00100 00010 (bin)

GeoHash results by grouping the sequence of bits by 5 and representing each group in base 32

10	11
00	01

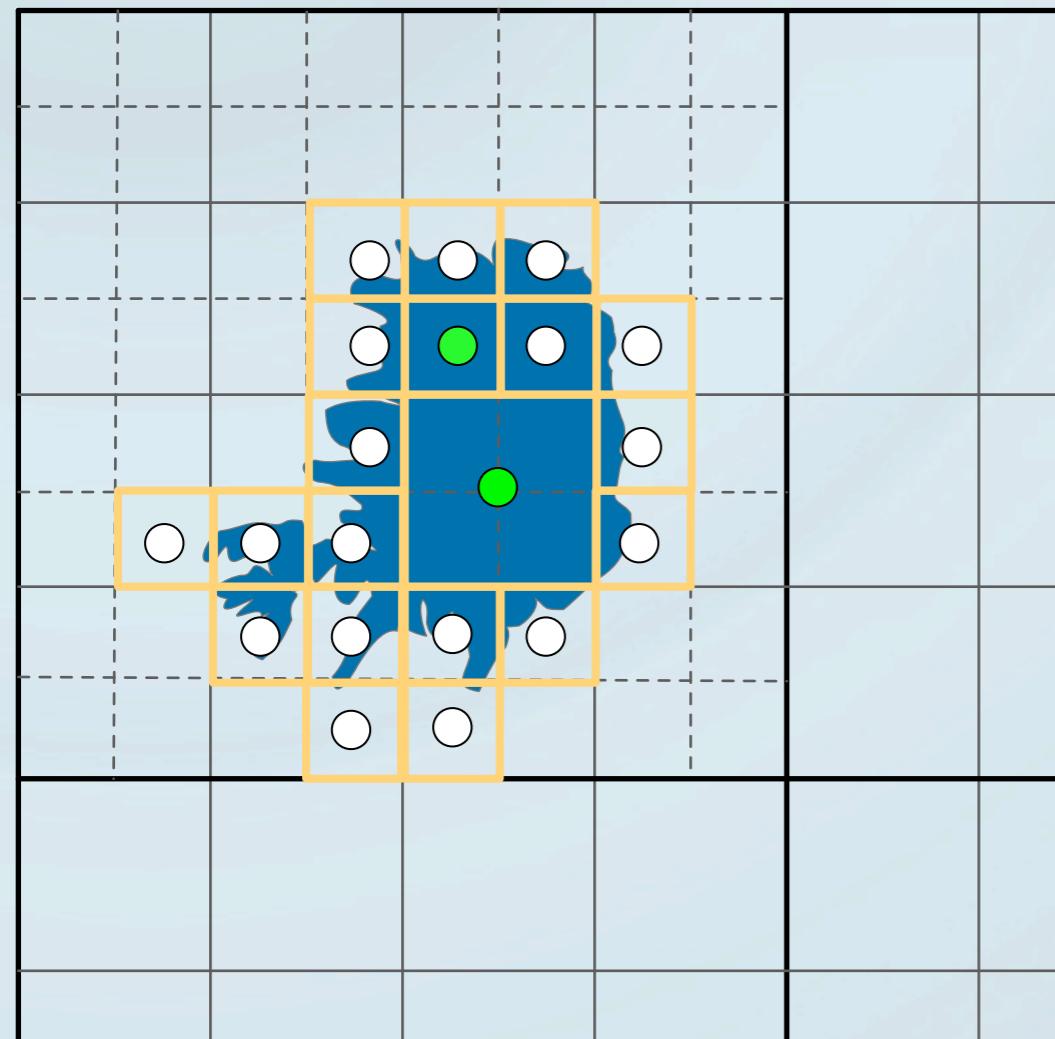
13 ₁₀ 31 ₁₀ 24 ₁₀ 4 ₁₀ 2 ₁₁₁₀ (dec) ₁₁₁ ⇒ ezs42			
10	00	11	01
1000	1001	1100	1101
0010	0011	0110	0111
0000	0001	0100	0101

101010	101011	101110	101111			
101000	101001	101100	101101			
100010	100011	100110	100111			
100000	100001	100100	100101			

CHARACTERISTICS

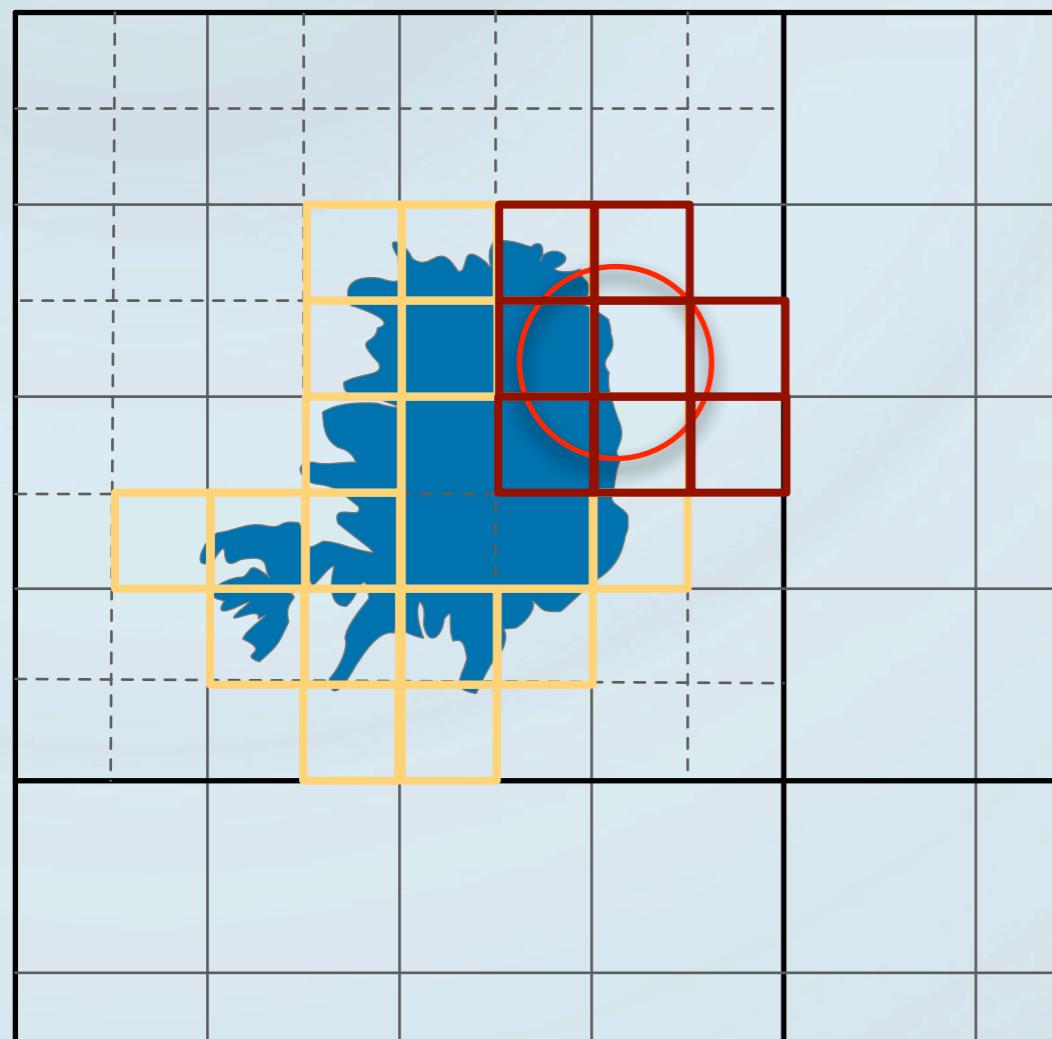
Any level of precision can be obtained by increasing the number of bits
Resolution can be dialed down or up by ignoring or not ignoring additional bits

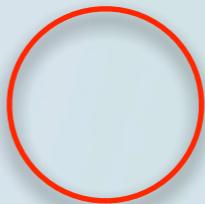
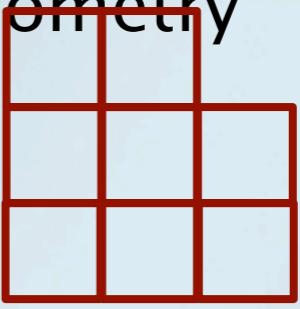
INDEXING A FEATURE



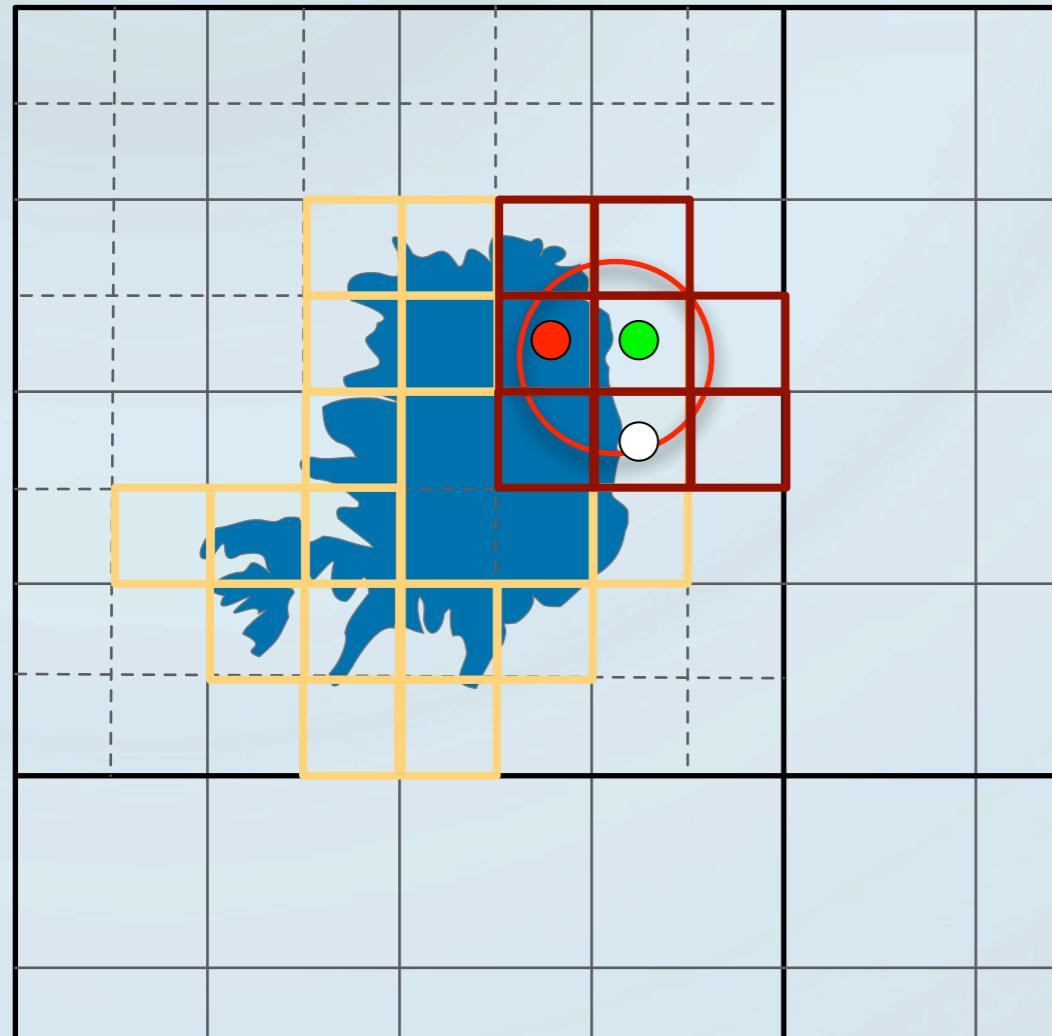
- ❑ If the feature is a point, simply index it at the highest resolution GeoHash (GH) that contains it
- ❑ Otherwise, find the GHs fully covering the feature (10)
- ❑ For each such GH, evaluate its children's relationship with the feature geometry recursively
 - ❑ If disjoint, return
 - ❑ If fully containing, index the feature at this GH, noted as "inside"
 - ❑ If at the max depth below the start GH, index the feature at this GH, noted as "related"

QUERYING GEOMETRY



- Specify the query geometry

- Find GeoHashes completely inside or intersecting the query geometry

- Result is the set of query relevant GHs

QUERYING GEOMETRY



- ❑ For each query-relevant GH
 - ❑ If the query GH is inside the query geometry, return all features in this GH and all those of all its children
 - ❑ Otherwise
 - ❑ Return features for which the GH is inside or is related and overlaps the query geometry
 - ❑ Until the highest resolution is reached, repeat
 - ❑ Until lowest resolution is reached, consider the parent GH*
 - ❑ Return features for which the GH is inside, or is related and overlaps the query geometry

Move it forward. Make it happen.



KEY POINTS

- While a given geometry may be indexed at many levels of resolution, each portion of it is indexed at only 1 level of resolution
- Indexing scheme naturally floats geometries into smaller GHs (i.e. higher resolution)
 - Biggest GH for a geometry is on the scale of the geometry itself (the GH that fits inside the geometry)
 - Only extremely large geometries are indexed under very big GeoHashes
 - As you walk up in GH size (down in resolution) fewer geometries are returned, thus fewer overlap calculations need to be made
- Although many query-relevant GHs may intersect a given geometry (e.g. along its edges), only 1 overlap calculation is ever made for that geometry
 - Once we determined the geometry is relevant to the query, we don't have to check it for every query-relevant GH
- *For each query-relevant GH we explore all the child hashes, but only about 1/4 of them must explore all the parent hashes
- Indexes are sharded across the cluster and all index operations (including overlap calculation) are performed server side on all shards in parallel



GEOSPATIAL-TEMPORAL DATA & THE ATOM OF SEMAPP SEMANTICS



Move it forward. Make it happen.



GEOSPATIAL-TEMPORAL DATA

“EVERYTHING HAS A SOMEWHERE AND A SOME-WHEN”

- CDR J. SMITH, USN RETIRED

- Geometry represents geospatial information a graphic system
 - Three elementary figures of (planar) geometry: point, line, area
 - Represented by a delimited set of geographic coordinates (i.e. points)
- Feature
 - In a GIS, a feature is a geometry with a label and a type (e.g. 30°10'N 031°06'E, Nile, River)
 - In the SRF, a feature is a sign that also has a geometry (e.g. [Nile, River, 30°10'N 031°06'E])
 - Geospatial disambiguation
 - GIS features fit naturally into the sign-space
 - Sign-Space can easily emit features to GIS systems
 - Associations may be asserted between GIS features and other signs / associations in the sign-space
- Interval is the elementary figure representing temporal information
- Event
 - In the SRF, is a sign that has an interval (e.g. [Wall Street bombing, Terrorist Act, 19200916])
 - Temporal disambiguation
 - Events fit naturally into the sign-space
 - Sign-Space can easily emit events to ABI systems
 - Associations may be asserted between events and other signs in the sign-space
- Layers
 - In a GIS, a layer created by displaying a set of features by type (e.g. Rivers, Roads)
 - In SRF, a layer can include signs of any type and types from any domain
 - Imagine temporal layers!

GIS SYSTEM'S
SET OF TYPES COMPRIZE A
DATA-MODEL

Move it forward. Make it happen.

THING 5 IS A UNIVERSAL STORE FOR GEOSPATIAL TEMPORAL INFORMATION

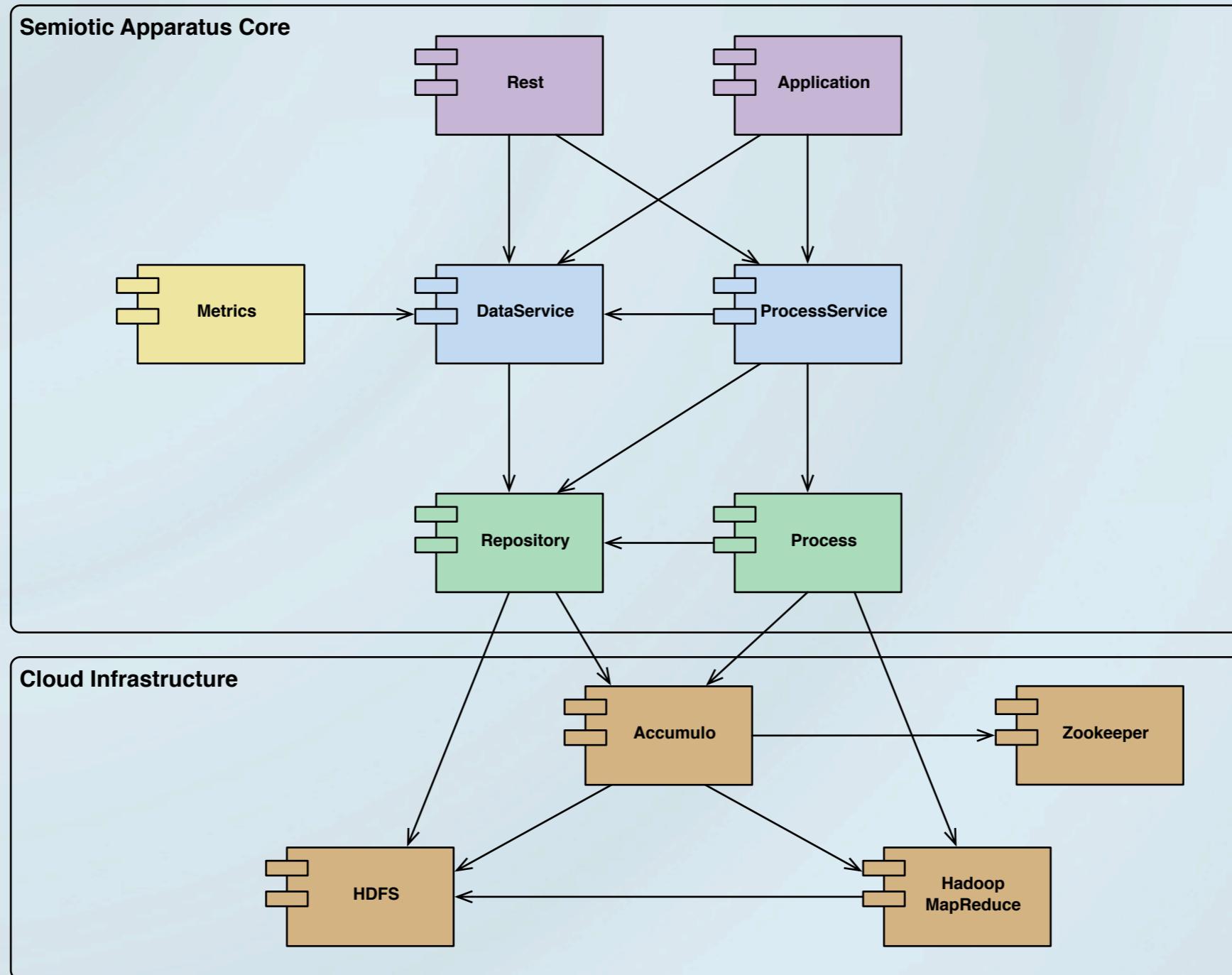


INFRASTRUCTURE



Move it forward. Make it happen.

LOGICAL ARCHITECTURE CORE

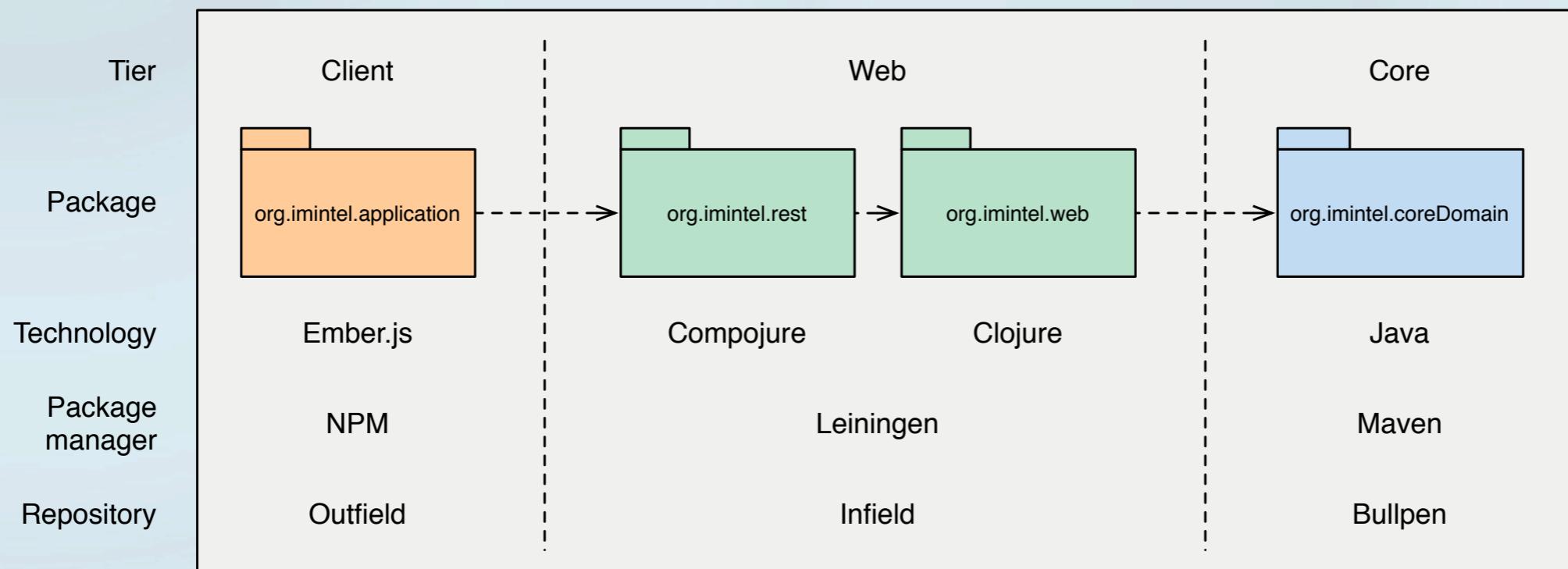


SemApp Logical Architecture

R 0.14 2011 01 23
S. Yoakum-Stover



LOGICAL ARCHITECTURE WEB

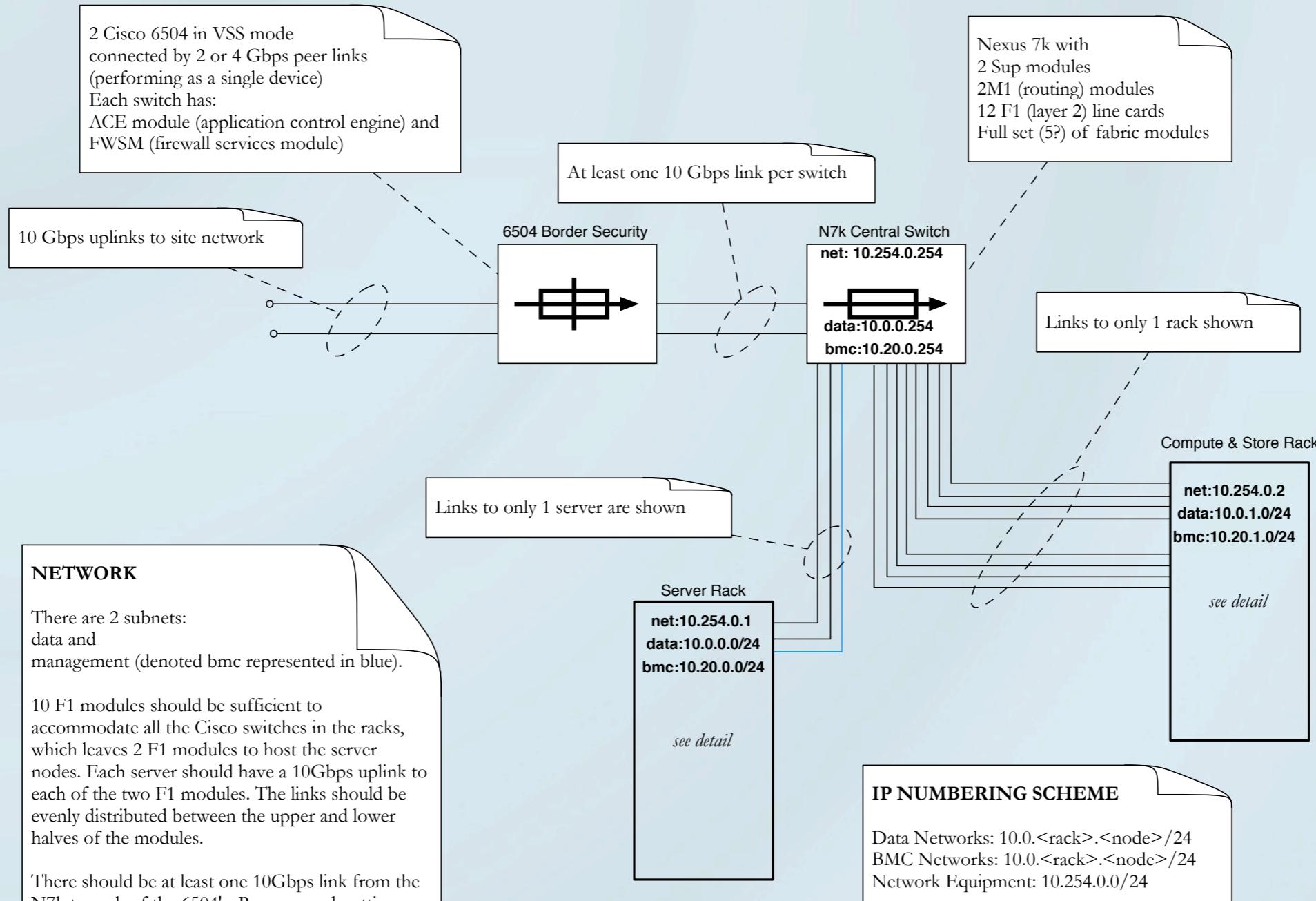


SemApp Web Architecture

R 0.17
S. Yoakum-Stover
2013 05 10

Move it forward. Make it happen.

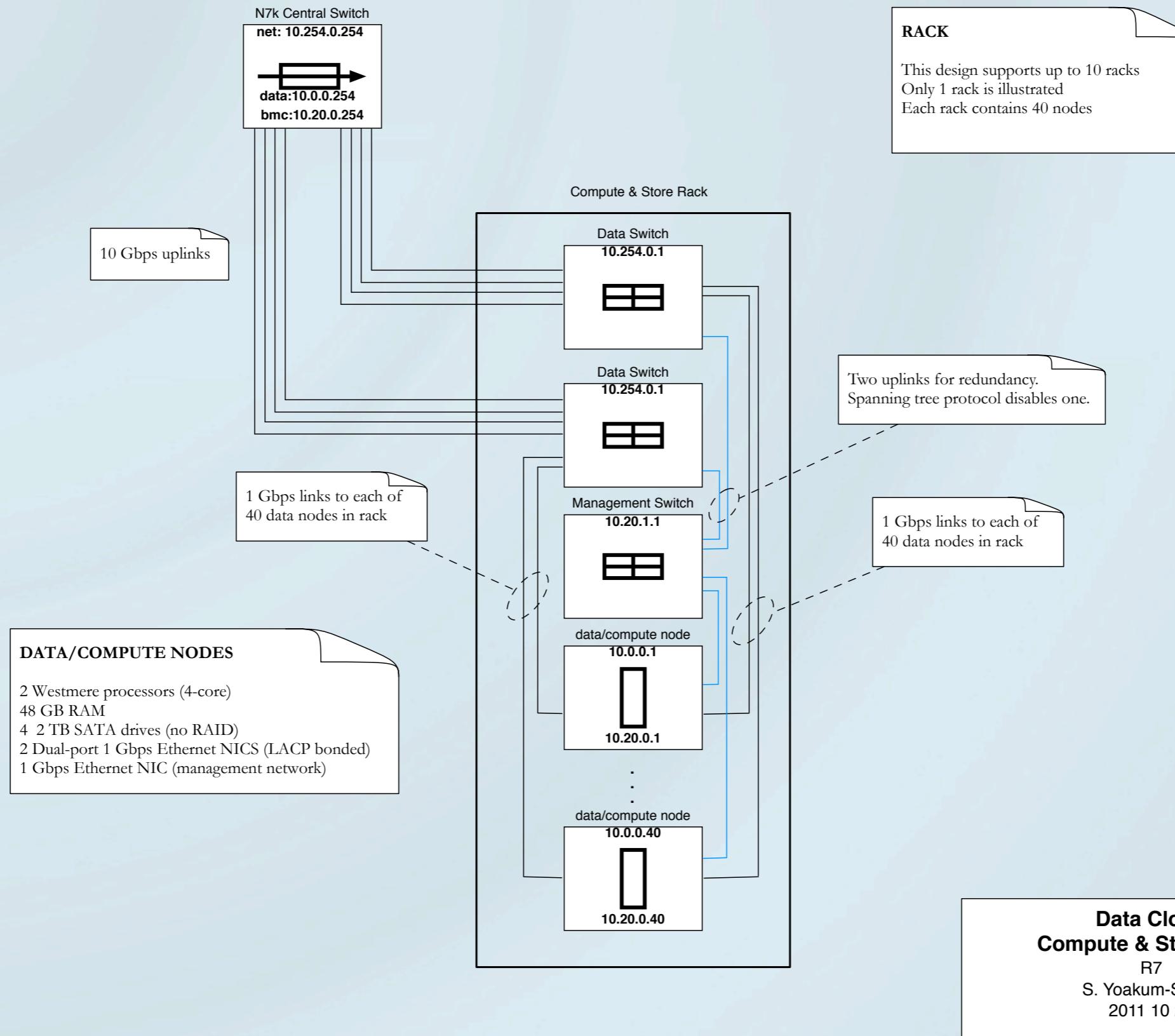
CENTRAL SWITCH



Design follows NSA
R6 (Not a Ghost
Machine)

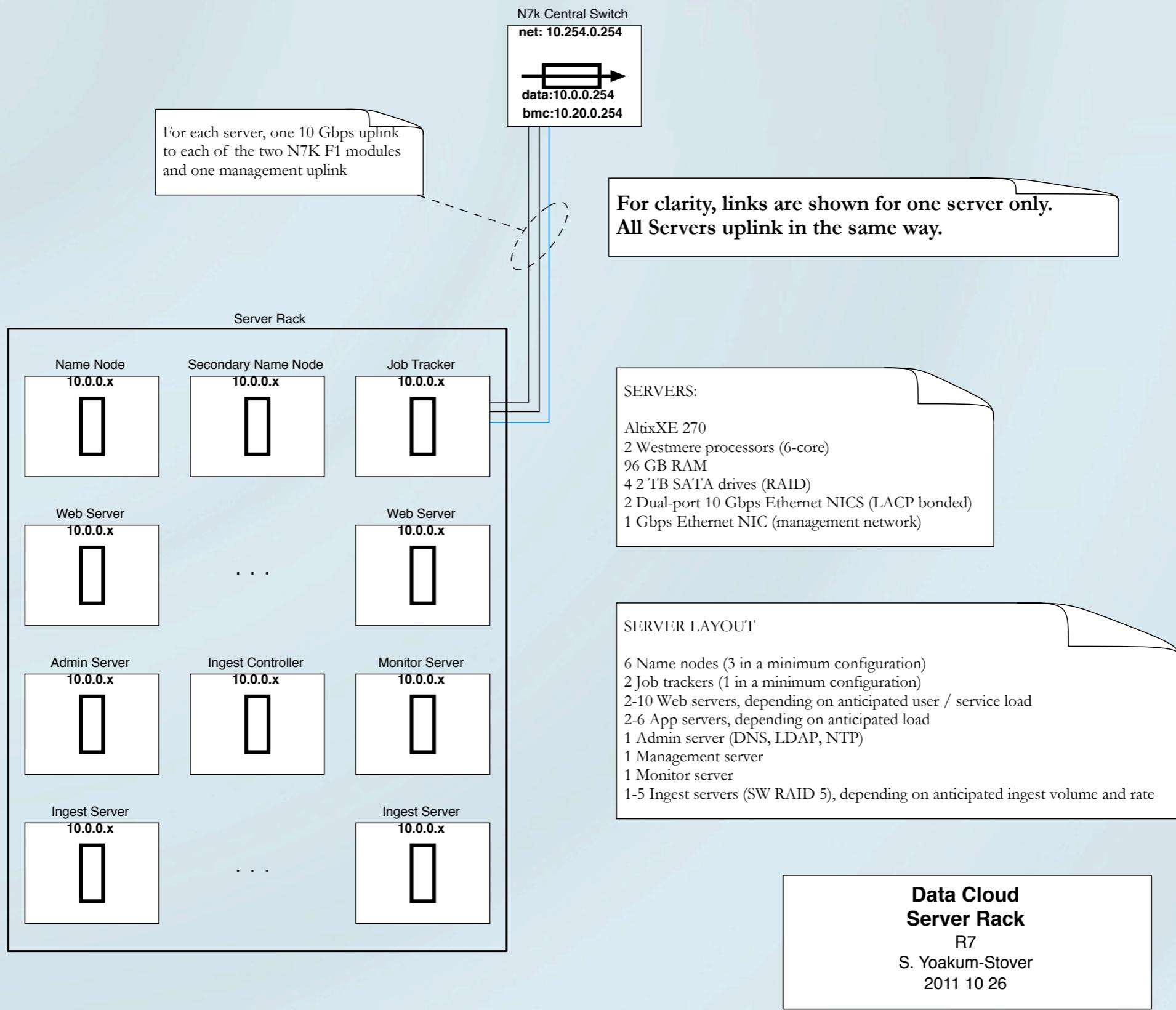
Data Cloud
Border Security and Central Switch
R7
S. Yoakum-Stover
2011 10 26

COMPUTE & STORAGE RACK





SERVER RACK





INGEST



Move it forward. Make it happen.

INGEST PROCESS

CORE TENANTS

- Never drop data – ingest is transactional
- Make it extremely robust – ingest it, then fix it
- As simple as possible and no simpler
 - Take the data as it lies – no transformation, no normalization
 - Move the bytes
 - Read bytes from endpoint (filesystem, JMS, (S)FTP, multi-cast...)
 - Persist raw bytes artifact
 - Persist a sign representing the artifact
 - Attach out-of-band metadata (ingest metadata)
 - Defined by the ingest manager and applied to all artifacts from a given endpoint
 - Fields
 - Required: authoritative source, ingest source, modality, classification and releasability
 - Optional: from-file-path, custom fields, detailed class info (e.g. system, declass, downgrade)
 - Persisted as signs associated to the artifact sign
- Everything else is a map-reduce process

Move it forward. Make it happen.



SOME TYPES OF INGEST

ARTIFACT INGEST

- ❑ Unstructured / semi-structured files of any type / modality
- ❑ NITF, IRC, email, MP2TS, ...
- ❑ Persist the artifact raw bytes
- ❑ Persist a sign representing the artifact
- ❑ Associate out-of-band metadata

RELATED ARTIFACT INGEST (E.G. SHAPE FILES)

- ❑ Shapes, index, shape attributes, and a decoder
 - Index relates shapes to attributes
 - Decoder relates attribute key-value pairs to signifier-concept pairs
- ❑ Ingest tarball of shapes, index, attributes and decoder as an artifact
- ❑ A map-reduce job persists shapes and attributes as signs and associations
- ❑ Spreadsheet ingest is similar

COLLECTION INGEST (E.G. WAMI)

- ❑ Maintains relationships among ingested artifacts

WEB SERVICE INGEST

- ❑ Collect data and model (service spider)
- ❑ Ingest the model, ingest the data as an artifact
- ❑ Process with map-reduce to reveal semantics
- ❑ Can be more sophisticated, incremental, dynamic

RDB INGEST

- ❑ Operates at the level of relational semantics (not domain semantics)
- ❑ Ingest the model
- ❑ Ingest the data and relations (i.e. signs and associations)

MODEL INGEST

- ❑ Represent model as OWL ontology
- ❑ Ingest the model
- ❑ Model, concepts, predicates are persisted with expressivity equal to OWL

Move it forward. Make it happen.



ORGANIZED TESTING



Move it forward. Make it happen.

TESTING

ORGANIZATION OF TESTS

- Core service boundary
 - Where the system data architecture is fully enforced
 - Dividing line between managed & non-managed tests
- Non- Managed tests
 - Unit tests
 - Result from TDD
 - Not included in test plans
- Managed tests
 - Included in test plans
 - Functional
 - Configuration
 - UI
 - Performance
 - Acceptance
- Regression testing
 - All automated tests are part of continuous integration process
 - Check-in – build – deploy – test – manual test

Type	Visibility	Mode	Level	Subject
Unit	white box	automated	repository process	CRUD SemApp elements
Functional	black box	mostly automated	service	data services,
Configuration	black box	various	n/a	infrastructure
UI	black box	various	application	UI
Performance	black box	various	various	user interactions,
Acceptance	black box	manual	application	applications



PARALLEL PROCESS VIDEO



Move it forward. Make it happen.



PARALLEL PROCESS VIDEO?

WHY YES!

Ruminate

Segment out the metadata groups

Segment out the frames

Expose the structural elements so they may be more easily exploited

Digest

Extract signs and associations representing the static metadata

The first group is all you need

Now it's discoverable

Process

Extract the collection footprint

Now it's discoverable from a geospatial query and we can do collection management

Exploit the pixels (e.g. object extraction)

etc...

Move it forward. Make it happen.



PIPELINES VS WAVES



Move it forward. Make it happen.

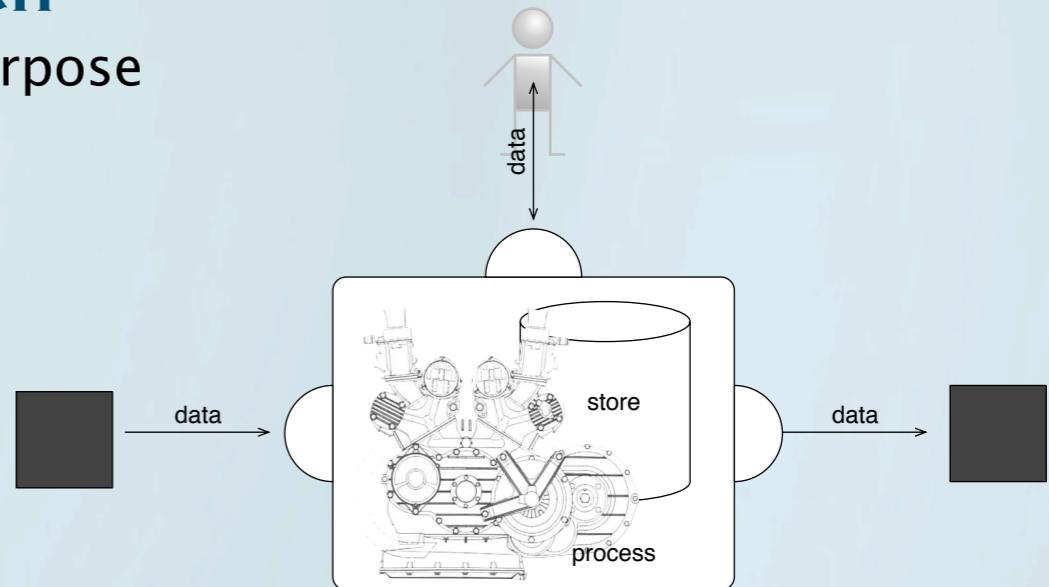
PROCESSING PIPELINES

DIVERSITY OF INFORMATION SYSTEMS, EACH

- Has its own processing, data representation, purpose
- Gets data from somewhere
- Transforms / creates data
- Outputs data

PROCESSING PIPELINE

- Common pattern used within systems
 - Used to chain modular process components
 - Examples: NLP, ETL-style ingest system
 - Data flows down the pipe
 - At every point in the pipe, the data is different
- Common pattern used in Enterprise integration
 - Services hooked up to fulfill business processes (SOA, BPEL, ESB)
 - Data flows between systems



ISSUES

- If processing on a piece of data fails
 - Subsequent processing can not occur (bad data)
 - Dropped on the floor / handled on dead letter channel
- If a component fails
 - Processing can not continue, data stops flowing (broken pipe)
- If a new component is added, downstream data is not affected
- The data assets are different at every point in the Enterprise
 - Collecting into separate pools of information
- Bandwidth & Complexity (IT, data management)

PROCESSING PIPELINES

DIVERSITY OF INFORMATION SYSTEMS, EACH

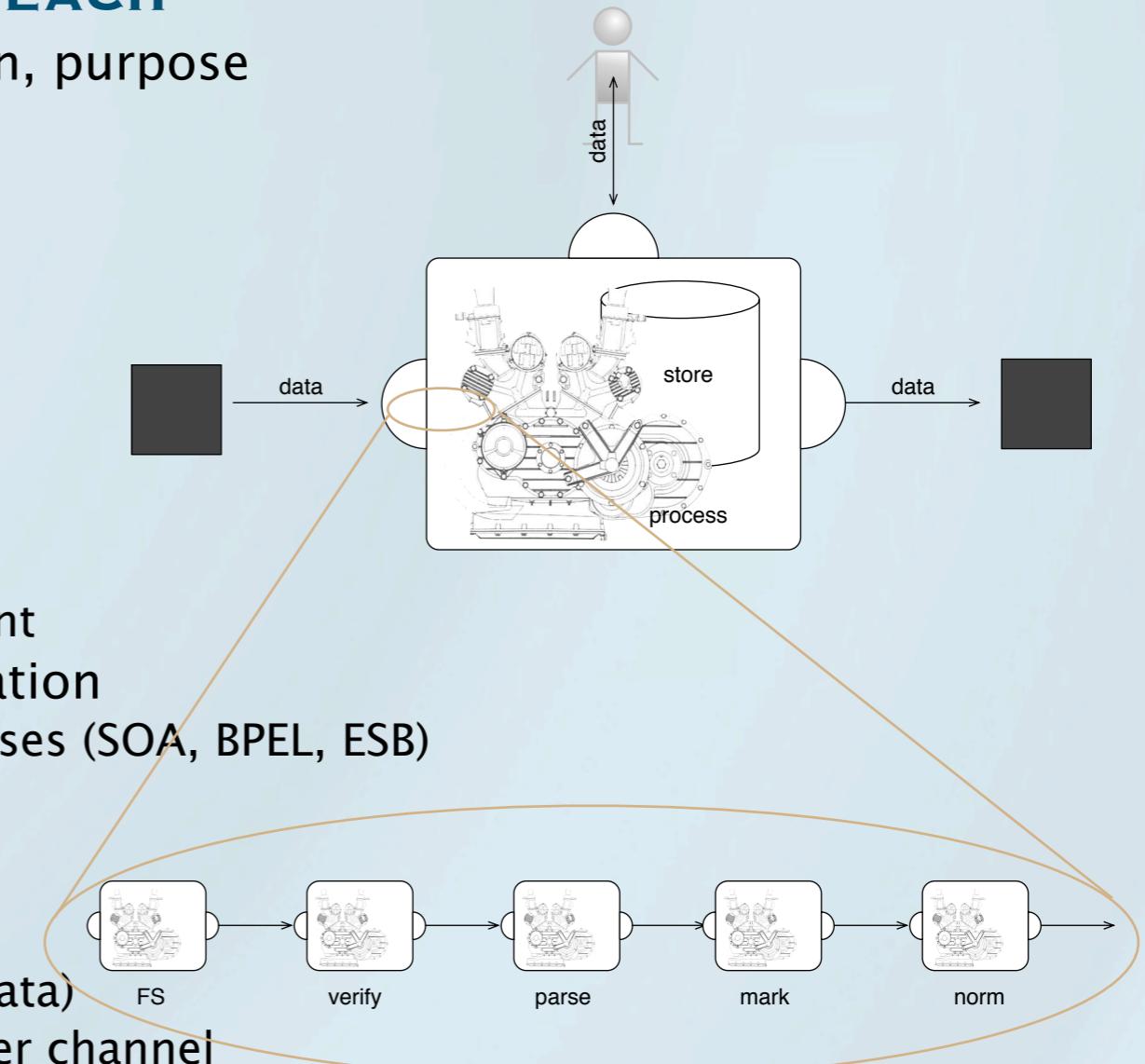
- Has its own processing, data representation, purpose
- Gets data from somewhere
- Transforms / creates data
- Outputs data

PROCESSING PIPELINE

- Common pattern used within systems
 - Used to chain modular process components
 - Examples: NLP, ETL-style ingest system
 - Data flows down the pipe
 - At every point in the pipe, the data is different
- Common pattern used in Enterprise integration
 - Services hooked up to fulfill business processes (SOA, BPEL, ESB)
 - Data flows between systems

ISSUES

- If processing on a piece of data fails
 - Subsequent processing can not occur (bad data)
 - Dropped on the floor / handled on dead letter channel
- If a component fails
 - Processing can not continue, data stops flowing (broken pipe)
- If a new component is added, downstream data is not affected
- The data assets are different at every point in the Enterprise
 - Collecting into separate pools of information
- Bandwidth & Complexity (IT, data management)



PROCESSING PIPELINES

DIVERSITY OF INFORMATION SYSTEMS, EACH

- Has its own processing, data representation, purpose
- Gets data from somewhere
- Transforms / creates data
- Outputs data

PROCESSING PIPELINE

- Common pattern used within systems
 - Used to chain modular process components
 - Examples: NLP, ETL-style ingest system
 - Data flows down the pipe
 - At every point in the pipe, the data is different
- Common pattern used in Enterprise integration
 - Services hooked up to fulfill business processes (SOA, ESB)
 - Data flows between systems

ISSUES

- If processing on a piece of data fails
 - Subsequent processing can not occur (bad data)
 - Dropped on the floor / handled on dead letter channel
- If a component fails
 - Processing can not continue, data stops flowing (broken pipe)
- If a new component is added, downstream data is not updated
- The data assets are different at every point in the End-to-End pipeline
 - Collecting into separate pools of information
- Bandwidth & Complexity (IT, data management)





THE PROBLEM WITH PROCESSING

is...

We got **Lots** of Different Processes and **Lots** of Data to process

and Everyone **Wants** everyoneElse's processed Data

So we we Hook Up the processes and Send the Data

In other words,

The Data **Flows** around the Enterprise
(or NOT)

and is Different at Every Place

So **Don't** do That



PROCESSING WAVES

SEMAPP DECOUPLES PROCESSING FROM DATA FLOW

- ❑ Ingest is Not a pipeline
 - A single process: Move file/stream bytes into the sign-space (and persist OOB metadata)
 - Works the same regardless of what sort of thing the bytes are
 - Everything else is a process
- ❑ Sign-Space is an ocean of data
 - Data stays in place
 - Processing enriches the ocean
 - All processing sees the same ocean

PROCESSING OCCURS IN WAVES

- ❑ Each wave processes all the data available to it
 - Available means: Is processable & Not already processed
- ❑ Processes leave their fingerprints on the things processed
 - Mechanism for determining availability during the next wave / for another process
- ❑ Processes execute independent of other processes

NO FLOW, NO LEAKS, NO BREAKS

- ❑ All processing that can occur on an element will
- ❑ The sign-space is the floor, data is never dropped (i.e. lost)
- ❑ If a component fails, dependent processing finds no available input
 - Once fixed, dependent processing again finds available input on the next wave
- ❑ If processing on a piece of data fails, dependent processing will not find it to be available
 - Enrichment produced by all previous processing still present
 - If the issue can be corrected, no special handling is required, dependent processing find it on the next wave
- ❑ If a new process is added, all the data that can be processed will
- ❑ Bandwidth reduced, complexity simplified, and robust

Move it forward. Make it happen.



TEMPORAL INDEXING



Move it forward. Make it happen.

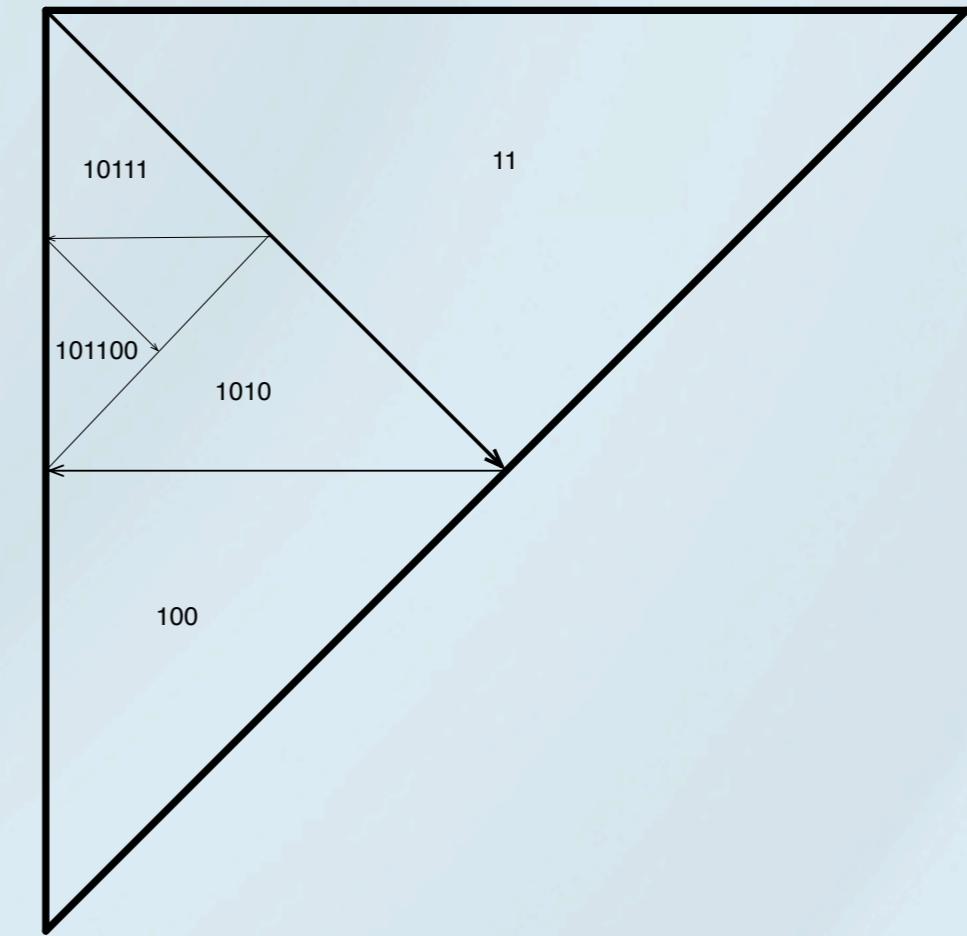
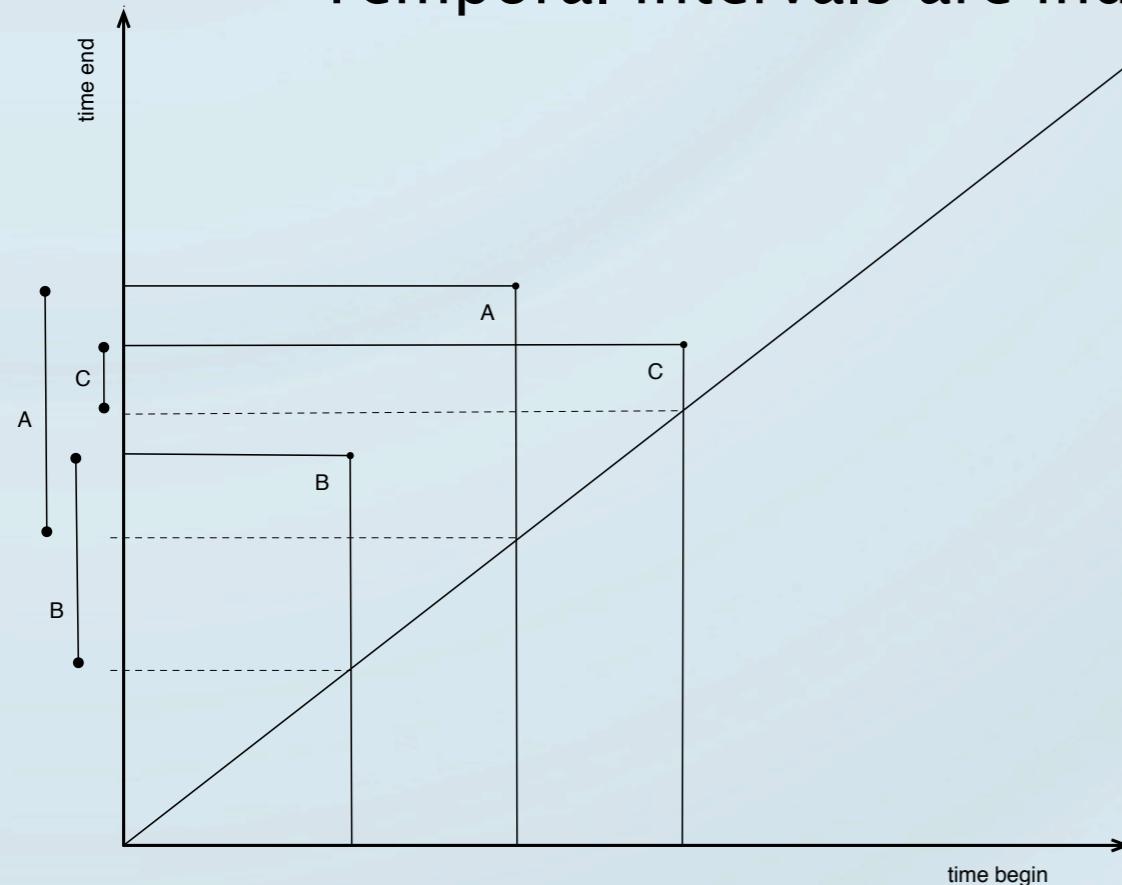
TD-TREE OVERVIEW

TEMPORAL DECOMPOSITION TREE* (TDT) REPRESENTS TEMPORAL INTERVALS AS POINTS IN 2D

Region of interest is a triangle, which is decomposed into a tree of ever smaller triangles

Each represented with a binary identifier

Temporal intervals are indexed in the leaves of the tree



*Stantic, B., Topor, R., Terry, J., Sattar, A.: Advanced Indexing Technique for Temporal Data. Computer Science and Information Systems, Vol. 7, No. 4, 679-703. (2010)

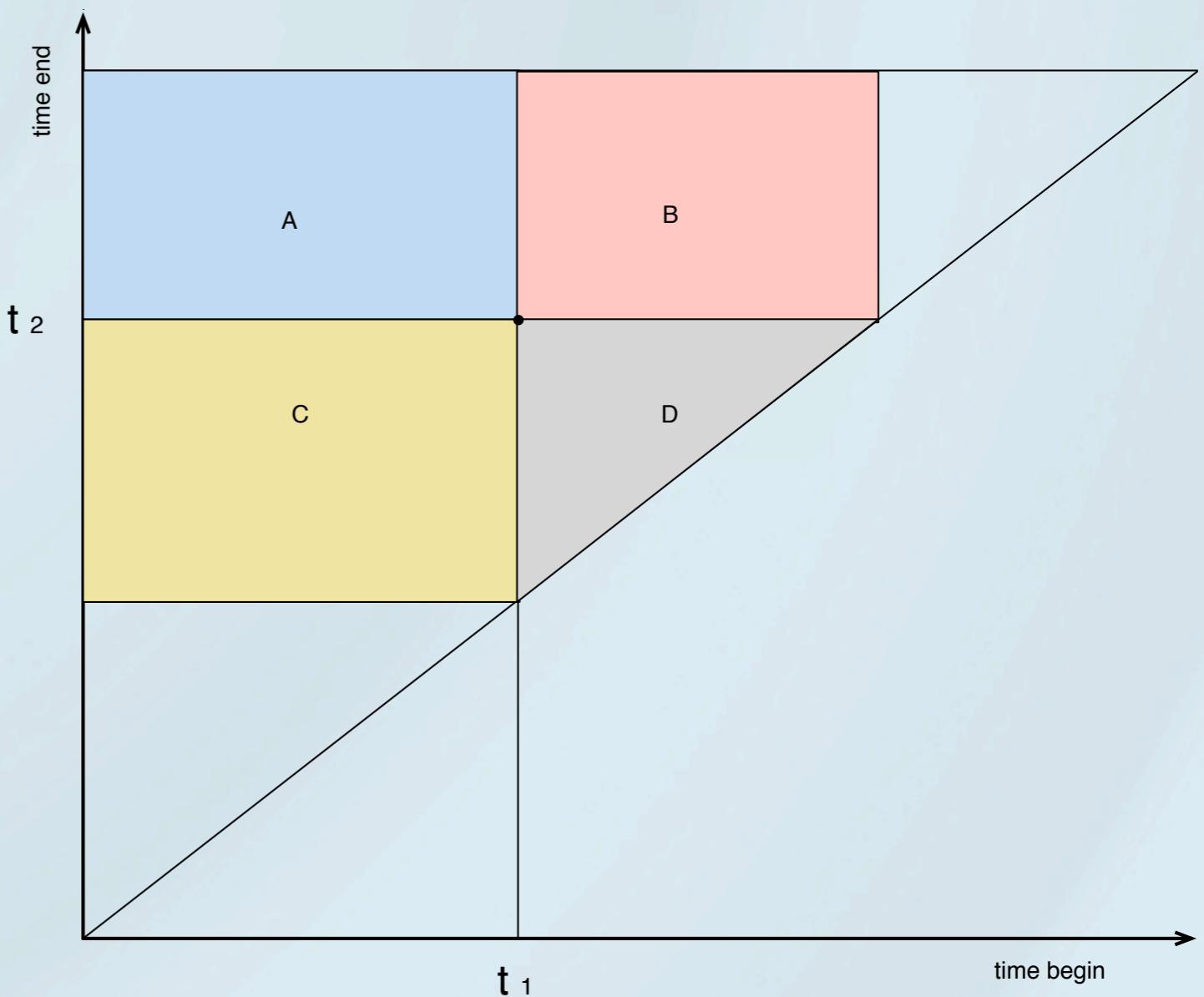
TEMPORAL INDEXING & QUERYING

INDEXING

- Index the interval at the TDT leaf which contains it (our approach differs slightly from the TDT paper)

QUERYING

- Standard queries expressed as rectangles, for example:
 - A: Intervals starting before t_1 and ending after t_2
 - B: Intervals starting between t_1 and t_2 and ending after t_2
 - C: Intervals starting before t_1 and ending between t_1 and t_2
 - D: Intervals starting between t_1 and t_2 and ending between t_1 and t_2





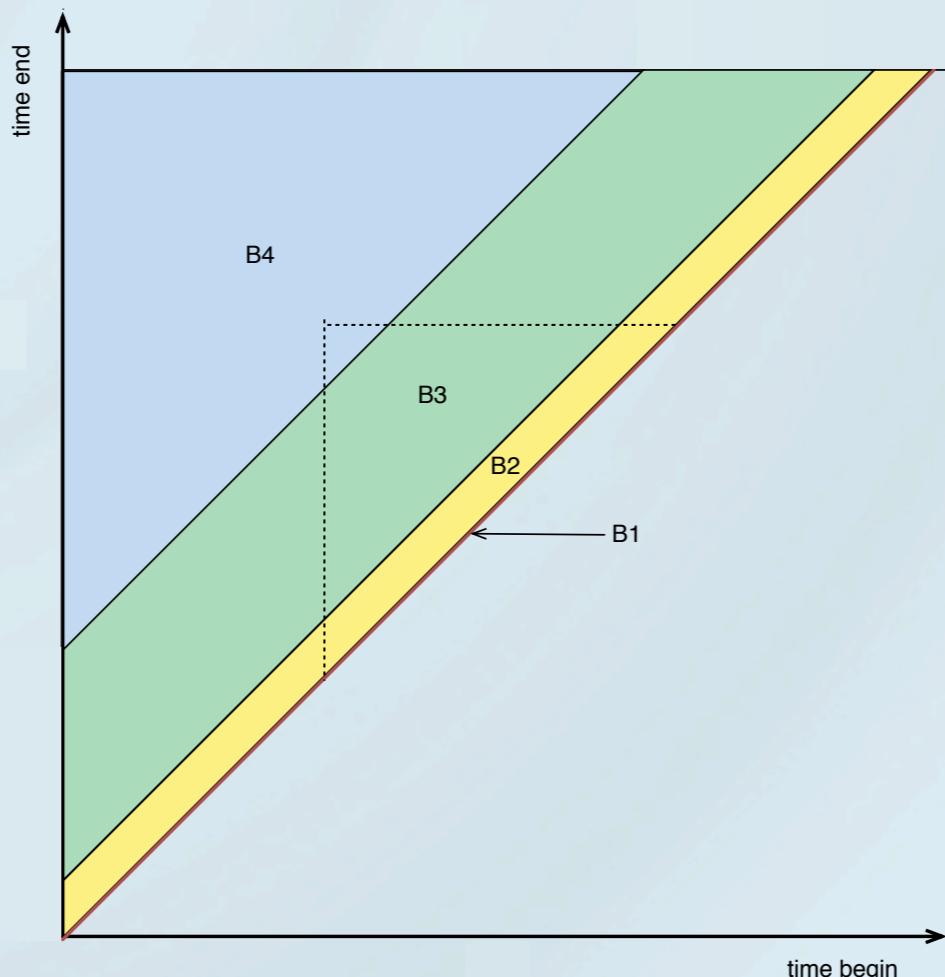
KEY POINTS

- Earliest time, latest time, and max resolution are specified prior to constructing the TDT, however
 - A large offset is used in case we need to go back farther in time
 - Max resolution can be increased later on
- Indexes are sharded across the cluster and all index operations are performed server side on all shards in parallel
 - Keeps the index balanced across the cluster
- TDT approach works best for intervals that are widely distributed in size
 - Points in time are more efficiently indexed using straight-forward approach
 - Intervals relevant to Intel will likely cluster just above the $t_1 = t_2$ line
 - Few intervals are likely to appear in the future

Move it forward. Make it happen.

TEMPORAL BUCKET INDEX

- Addresses the shortcomings of the TDT index
- Separately index intervals with different lengths
 - Instants (intervals with 0 length) should just be lookups (B1)
 - Longer intervals may require some small amount of analysis (B2, B3, B4, ...)



- Index
 - Assert element in the smallest bucket size that fully contains its interval, ordered by start time
- Query
 - For each bucket, starting with the earliest start time that could satisfy our query and ending with the



THERE IS NO METADATA



Move it forward. Make it happen.



METADATA

- Metadata is structured data
 - Generally distinguished by its target
 - Not the X, but information about (or extracted from) the X, where X = artifact, data, model ...
 - Subjective distinction – One person's metadata is another's data
 - Variety of metadata standards/models reflect a variety of uses and perspectives
- SemApp accommodates all, but sanctifies none
 - Metadata is represented just like any other kind of structured data
 - Metadata and metadata models are persisted in SRF
 - SemApp imposes no special metadata model
- Pedigree, logs, and system actors are part of the SemApp machinery
 - History of sign-space changes (who, what, when) are persisted
 - System logs are persisted
 - System actor (person/process known to the SemApp) information is also persisted
 - Supports history, logging, and security

METADATA IS STRUCTURED DATA



TREATING SEMI-STRUCTURED DATA RIGHT



Move it forward. Make it happen.



SEMI-STRUCTURED DATA

Semi-structured data is data that conforms to a standard
(e.g. HTML, RSS, MIME, NITF, MP2 TS, ...)

Standard is a less glamorous cousin to data-model

Mix of

Semantic elements (e.g. NITF originating station Id, MIME in-reply-to)

Structural elements (e.g. NITF file header length)

How to handle it?

Preserve & expose both the structure & the semantics



BEYOND INGEST

Ingest

Stream raw bytes into the Sign-Space

Ruminate

Segment according to the standard

Typically only 1 way to do this for a given standard, but there are many standards

Exposes the structural elements so they may be more easily exploited

Digest

Extract signs and associations according to a semantic model based on the standard

Expose the semantic elements expressed by the standard

(create signs with mentions referring to appropriate segments)

May be done in slightly different ways

(different semantic models interpreting the standard)

Process

Extract semantics from unstructured elements based on a domain model (e.g. NLP)

Expose and enrich semantic elements extracted from unstructured elements

May be done in many different ways according to many different models



WHO ARE WE?



Move it forward. Make it happen.



WHO WE ARE

Institute for Modern Intelligence

IMI aims to change the world!
by developing
an Ultra-Large Scale systems infrastructure
for data-intensive computing & Data-intensive operations

Our mission:

Develop the Science, Practice, and Governance of Modern Intelligence

We are a non-profit, 501(c)3 organization
whose purpose is to fulfill our mission.

S. Yoakum-Stover & M. A. Eick are the IMI founders

Mission Focus

Mission Focus aims to build clean code that works for its customers.
We are an agile software development shop that takes
domain design and development as seriously as
system design and development.

Our tag line is:

Move it Forward, Make it Happen

We are a small, for-profit organization
whose purpose is to fulfill its customers' missions.
Our profits help to support the IMI.

M. A. Eick & S. Yoakum-Stover are CEO and Chief Scientist respectively

Core Domain

Mission Focus and IMI work mostly in the intelligence domain with DoD and IC customers / partners.
Our core domain addresses the storage, processing, and utilization of data in the context of immense volume and diversity.

We are experts in cloud computing and storage technology.
We invented the Sign Representation Framework which underpins a game-changing approach to data unification.

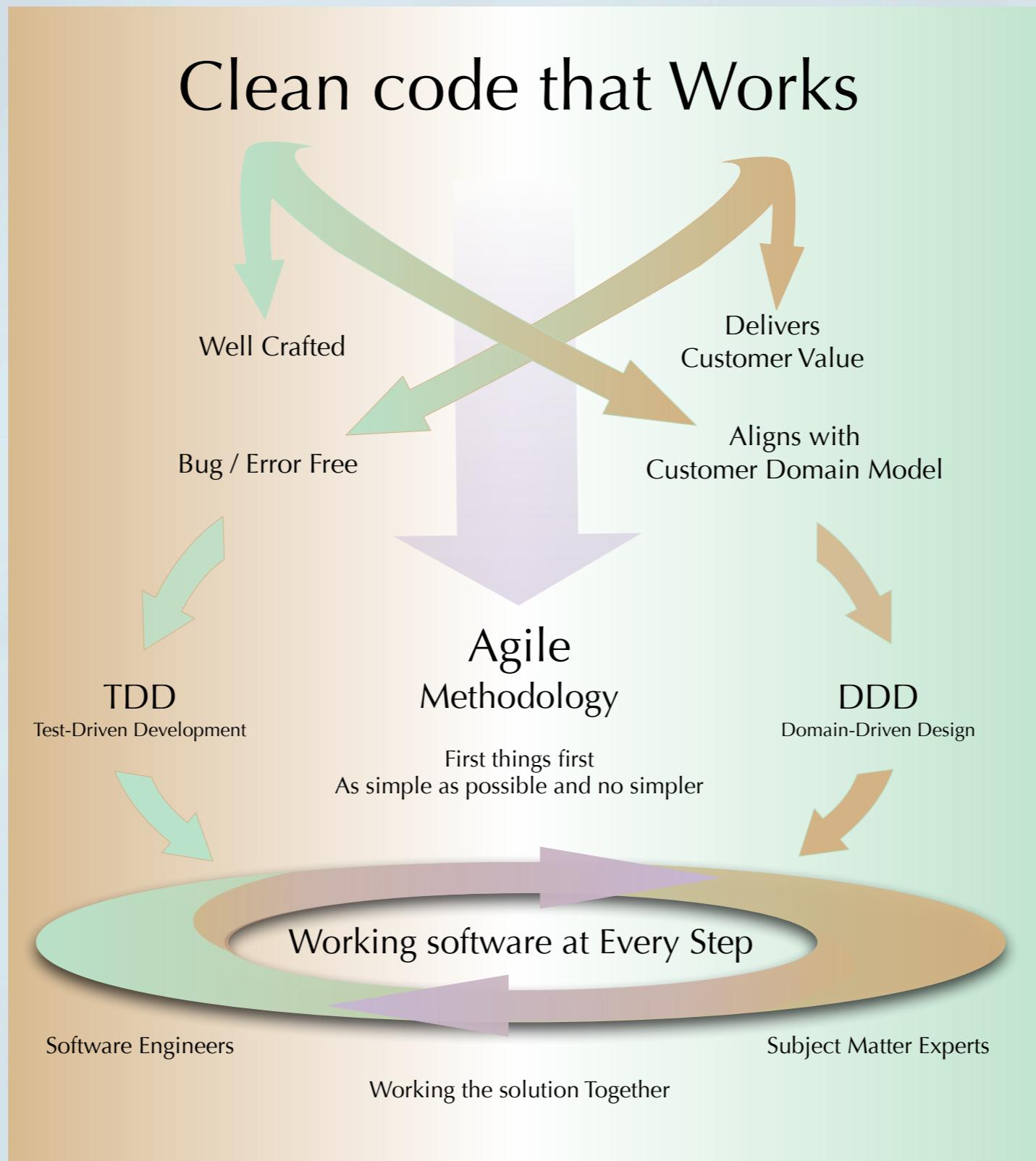
We pride ourselves on our disciplined engineering practices and distinguish ourselves by our ability to continually learn and innovate.

The work we do is meaningful and intentional and is wrapped with our integrity.

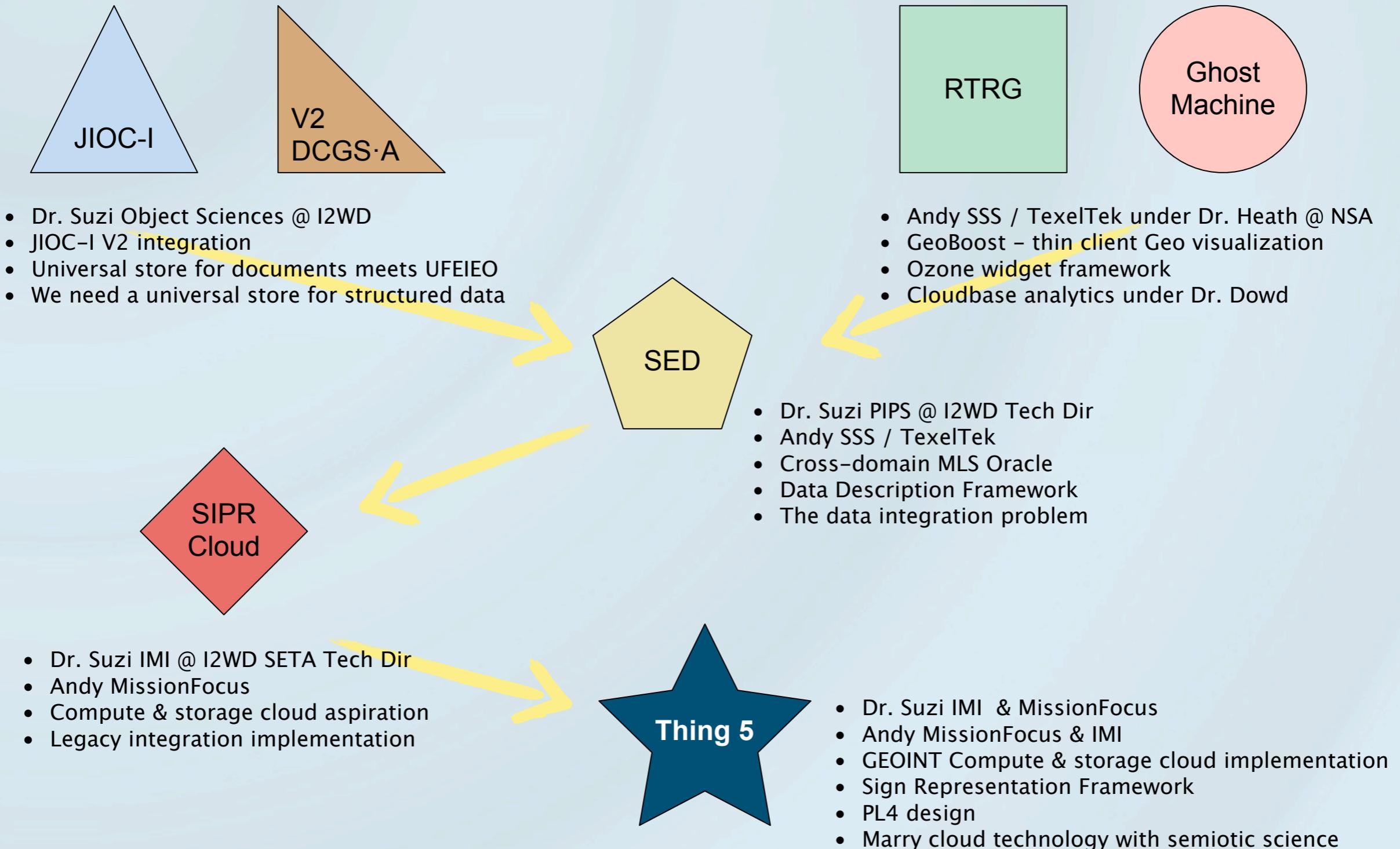
Put simply,
We just think harder and work better than the rest.



HOW WE WORK



THING 5 HERITAGE



BATTLE SCARRED AND A WHOLE LOT SMARTER

CURRENT CAPABILITIES



Move it forward. Make it happen.

CAPABILITIES - DATA

■ COLLECT, INGEST, EXPORT

- ALL UNSTRUCTURED / SEMI-STRUCTURED FILE TYPES
- FROM ENDPOINT TYPES: FILE SYSTEM, JMS TOPIC / QUEUE, (S)FTP, MULTI-CAST / UNI-CAST STREAM
- RELATED FILES (E.G. SHAPE)
- COLLECTIONS (E.G. WAMI)
- WEB SERVICE COLLECTORS
- DOMAIN DATA-MODELS AS OWL

■ ASSERT, RETRIEVE, DELETE, EXISTS, TALLY

- ARTIFACTS OF ANY TYPE AND MODALITY
- SIGNS & ASSOCIATIONS
- DOMAIN MODELS
- SYSTEM ACTORS
- PROCESS ELEMENTS

■ INDEX & QUERY

- FULL TEXT
- SEMANTIC (SIGNIFIER, CONCEPT, PREDICATE)
- GEOSPATIAL (INCLUDING COMPOUND GEOMETRIES)
- TEMPORAL

■ PRESERVE HISTORY

- ALL SIGN-SPACE CHANGES INCLUDING DELETES

■ SEMIOTIC REFINEMENT



CAPABILITIES - PROCESS

- **PROCESS WAVES INFRASTRUCTURE**
- **SEMAPP MAP-REDUCE PROCESS FRAMEWORK**
- **PARALLEL PROCESSES**
 - ASSESS
 - GATHER ARTIFACT, SRF, AND MODEL STATISTICS
 - DETECT
 - ARTIFACT STANDARD
 - ARTIFACT LANGUAGE
 - ENRICH
 - EXTRACT & EXPOSE PLACE NAMES FROM ARTIFACT TEXT
 - EXTRACT & EXPOSE SEMANTICS FROM EMAIL, NITF, IRC, TIFF, MP2TS, CSV
 - EXTRACT & EXPOSE NITF IMAGES AS ARTIFACTS
 - EXTRACT & EXPOSE EMAIL ATTACHMENTS AS ARTIFACTS
 - EXTRACT & APPLY SECURITY MARKINGS FROM NITF & MP2TS
 - EXTRACT & EXPOSE VIDEO COLLECTION FOOTPRINT FROM MP2TS, KML FLY-OVER DISPLAY
 - FIND SHORTEST PATH BETWEEN 2 SEMANTIC ELEMENTS
 - EXTRACT & EXPOSE SHAPE-FILE FEATURES AND ATTRIBUTES
 - EXTRACT AND EXPOSE FACES WITHIN IMAGES
 - GENERATE STANDARD IMAGE REPRESENTATIONS (THUMBNAIL, SMALL, MEDIUM, FULL)
 - MODEL REASONERS - FUNCTIONAL, PREDICATECHAIN, TRANSITIVE
 - RELATEDNESS WITHIN ARCHIVES
 - INDEX
 - ARTIFACTS, SIGNS, ASSOCIATIONS, MENTIONS, MODELS
 - LOAD
 - SYSTEMACTORS, OWL MODELS, PROCESS SPECIFICATIONS, GAZETTEERS
 - MEND
 - MAINTAIN INTEGRITY OF INDEX AND MATERIAL SIGNS
 - UTIL
 - PARSE > 50 DOCUMENT, AUDIO, VIDEO, & IMAGE FILE TYPES

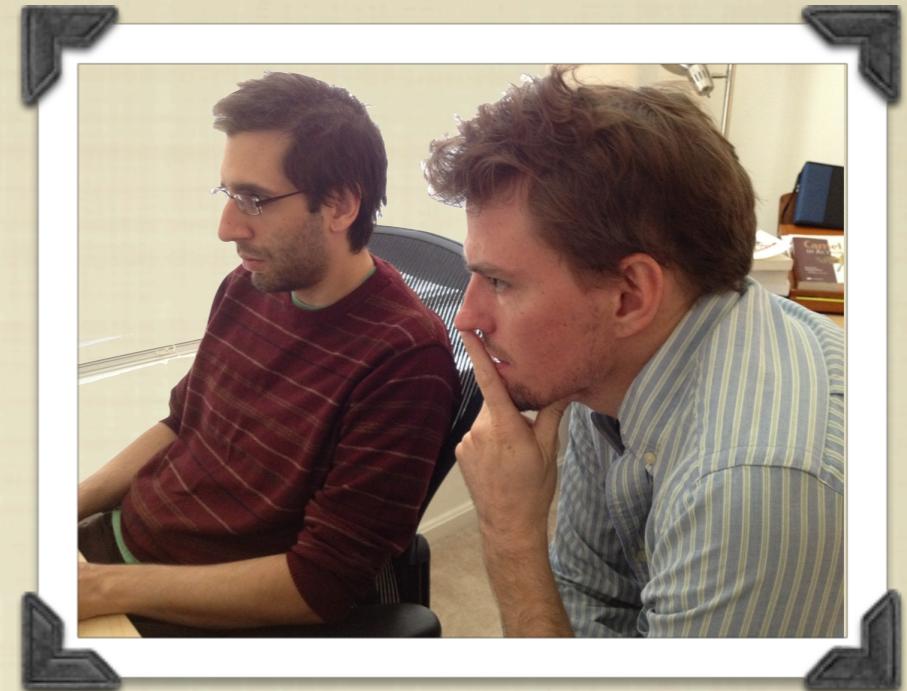




IMAGE & VIDEO PROCESSING

Move it forward. Make it happen.

Image Tag-Value Standard	Executable Process									
	Artifact Ingest	Artifact Standard Detector	Classification Markings Extractor	Artifact Ruminate	Artifact Digest	Artifact Text Extractor	Place Name Extractor	Image Proofer	Image Face Extractor	
KODAK_DCR	✓	✓		✓	✓	✓	✓			
KODAK_KDC	✓	✓		✓	✓	✓	✓	✓	✓	
EPSON	✓	✓		✓	✓	✓	✓	✓	✓	
HASSELBLAD	✓	✓		✓	✓	✓	✓	✓	✓	
NIKON	✓	✓		✓	✓	✓	✓	✓	✓	
OTHER_RAW	✓	✓		✓	✓	✓	✓			
BMP	✓	✓		✓	✓	✓	✓	✓	✓	
GIF	✓	✓		✓	✓	✓	✓	✓	✓	
JPEG	✓	✓		✓	✓	✓	✓	✓	✓	
PNG	✓	✓		✓	✓	✓	✓	✓	✓	
PSD	✓	✓		✓	✓	✓	✓			
TIFF	✓	✓		✓	✓	✓	✓	✓	✓	
NITF2_0	✓	✓	✓	✓	✓	✓	✓	✓	✓	
NITF2_1	✓	✓	✓	✓	✓	✓	✓	✓	✓	

Standard	Executable Process									
	Artifact Ingest	Artifact Standard Detector	Classification Markings Extractor	Artifact Ruminate	Artifact Digest	Artifact Text Extractor	Place Name Extractor	Footprint Creator	Video Proofer	
MP2TS	✓	✓	✓	✓	✓	✓	✓	✓	✓	



TEXT PROCESSING

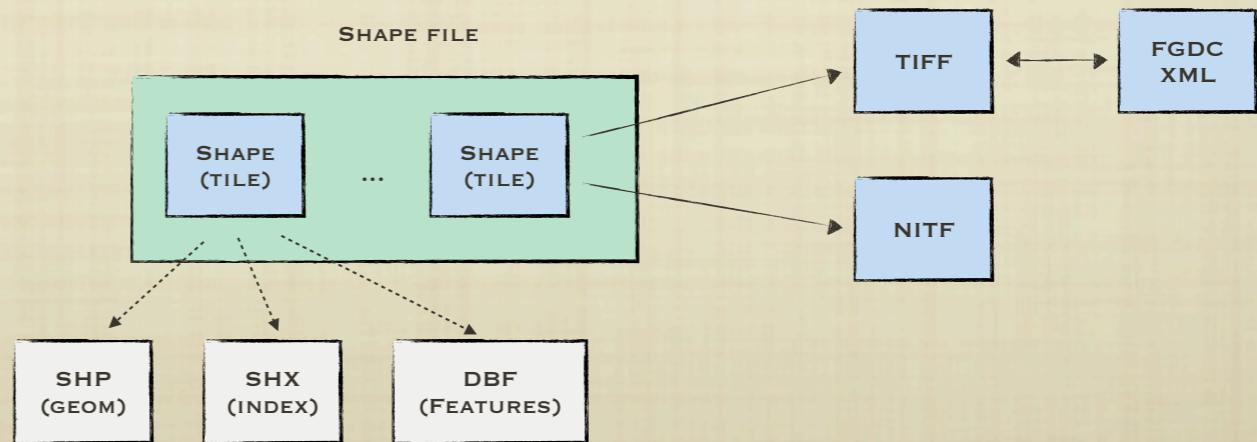
Move it forward. Make it happen.

Text Standard	Executable Process								FGDC Metadata Digest
	Artifact Ingest	Artifact Standard Detector	Artifact Ruminate	Artifact Digest	Artifact Text Extractor	Place Name Extractor	Artifact Language Detector		
CSV	✓	✓	✓	✓	✓	✓	✓	✓	
XLS	✓	✓	✓	✓	✓	✓	✓	✓	
ODS	✓	✓	✓	✓	✓	✓	✓	✓	
PLAIN_TEXT	✓	✓				✓	✓	✓	
XML	✓	✓				✓	✓	✓	
RSS	✓	✓				✓	✓	✓	
ODT	✓	✓				✓	✓	✓	
OTT	✓	✓				✓	✓	✓	
STW	✓	✓				✓	✓	✓	
SXW	✓	✓				✓	✓	✓	
FGDC_METADATA	✓	✓				✓	✓	✓	✓
HTML	✓	✓				✓	✓	✓	
DOC	✓	✓				✓	✓	✓	
DOCX	✓	✓				✓	✓	✓	
CHAT	✓	✓	✓	✓	✓	✓	✓	✓	
EMAIL	✓	✓	✓	✓	✓	✓	✓	✓	
RTF	✓	✓				✓	✓	✓	
PPT	✓	✓				✓	✓	✓	
PDF	✓	✓				✓	✓	✓	
PAGES	✓	✓				✓	✓	✓	
EPUB	✓	✓				✓	✓	✓	

ABOUT FGDC (FEDERAL GEOGRAPHIC DATA COMMITTEE)

LIDAR COLLECTION (SAIC) INCLUDES SHAPE FILES, TIFs, XML, AND NITFs.

EACH SHAPE FILES REFERENCES A TIFF WHICH REFERENCES AN FGDC XML FILE CONTAINING METADATA ABOUT THE SHAPE WHOSE ASSOCIATED IMAGE IS IN A NITF





PROCESSING (MISC.)

Move it forward. Make it happen.

Standard	Executable Process								
	Artifact Ingest	Artifact Standard Detector	Classification Markings Extractor	Artifact Ruminate	Artifact Digest	Artifact Text Extractor	Place Name Extractor	Footprint Creator	Archive Exploder
SHAPE		✓			✓				
SHP	✓	✓							
SHX	✓	✓							
DBF	✓	✓							
MODEL	✓	✓							
MP3	✓	✓							
MIDI	✓	✓							
MP4	✓	✓							
RIFF	✓	✓							
UNKNOWN	✓	✓							
DECODER	✓	✓							
ZIP	✓	✓						✓	
GZIP	✓	✓						✓	
BZIP2	✓	✓						✓	
TAR	✓	✓						✓	

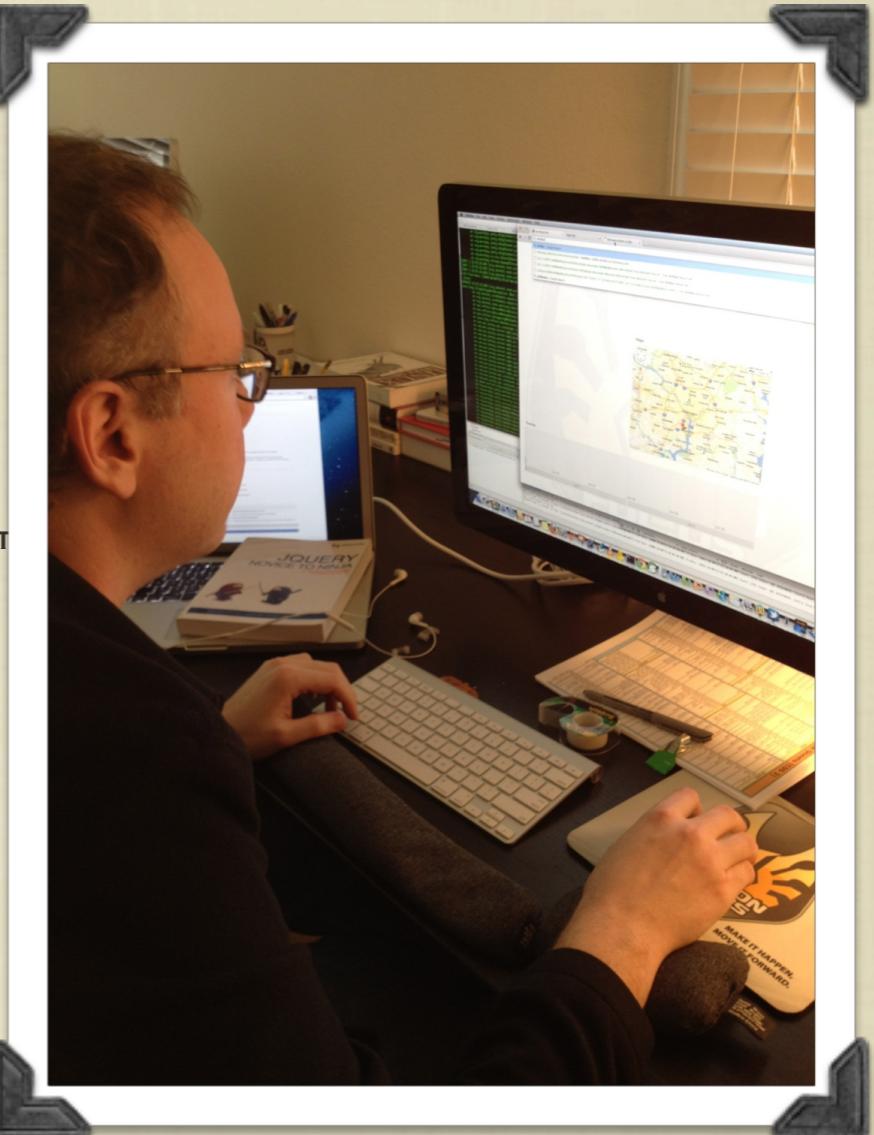
CAPABILITIES - APPLICATION

■ USER APPLICATIONS

- THE EYE (SEARCH APPLICATION)
- DATA SA (SURVEY APPLICATION)

■ PRIVILEGED USER APPLICATIONS

- INGEST MANAGER
 - CONFIGURE AND EXECUTE ARTIFACT INGEST
 - SPECIFY OUT-OF-BAND METADATA
 - SYSTEM ACTOR MANAGER
 - BROWSE & MANAGE USERS / ORGANIZATIONS / APPLICATIONS KNOWN TO THE SYSTEM
 - PROCESS MANAGER
 - SCHEDULE AND EXECUTE PROCESSES
 - BROWSE PROCESS ISSUES
 - MODEL MANAGER
 - GAZETTEER MANAGER
 - LAUNCHER
- API
- SEMAPP MAP/REDUCE FRAMEWORK
 - RESTFUL DATA & QUERY WEB SERVICES
 - KML EXPORT
- DOCUMENTATION
- DEPLOYMENT GUIDE, MANAGEMENT GUIDE, REST API



SEARCH UI

Text Search:

Text:

Search Examples

Will Return any data space element or artifact entailing this text.
Example: Suzi, Andy, Car

Semantic Search:

Text:

Search Examples

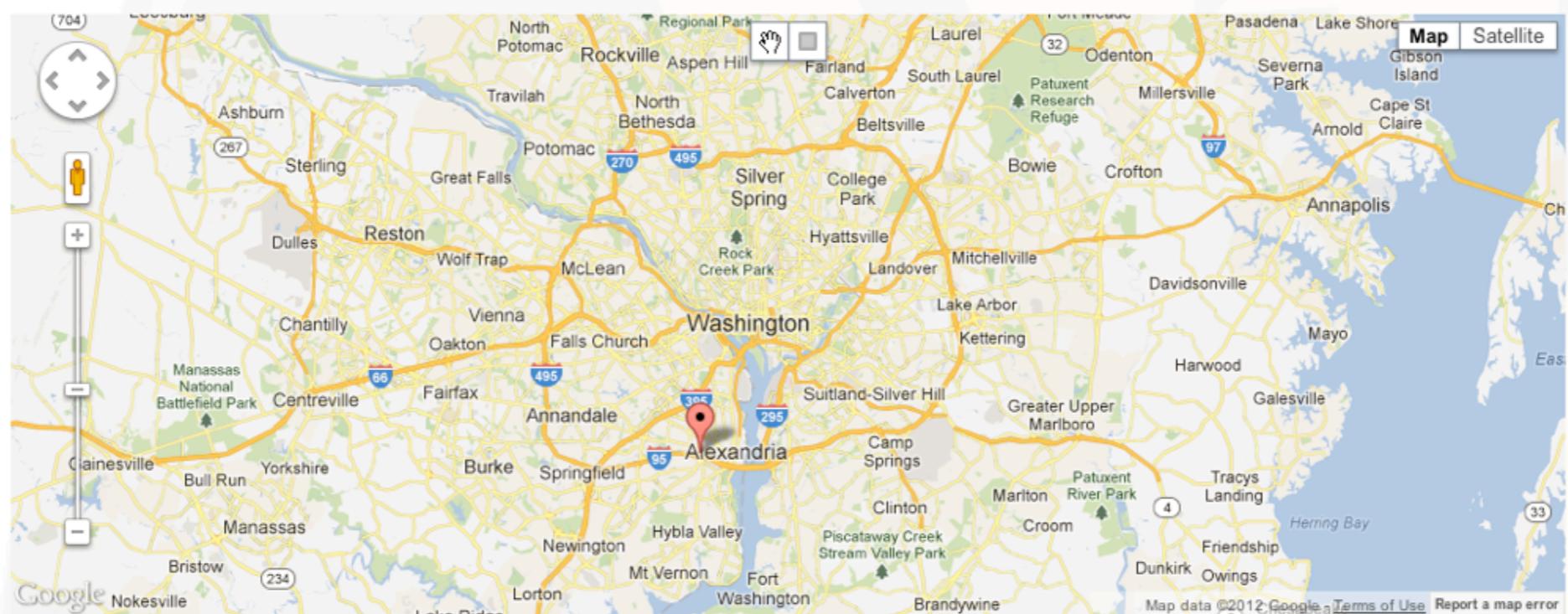
Will Return semantic elements that contain the following text in the sign. Example: Suzi, Andy, Artifact

Concept:

Will Return semantic elements that contain the following concept.
Example: Email, Sender, Receiver

Visual Search:

Select an area to Search



Temporal Search:

From

to

Search

RESULTS PAGE

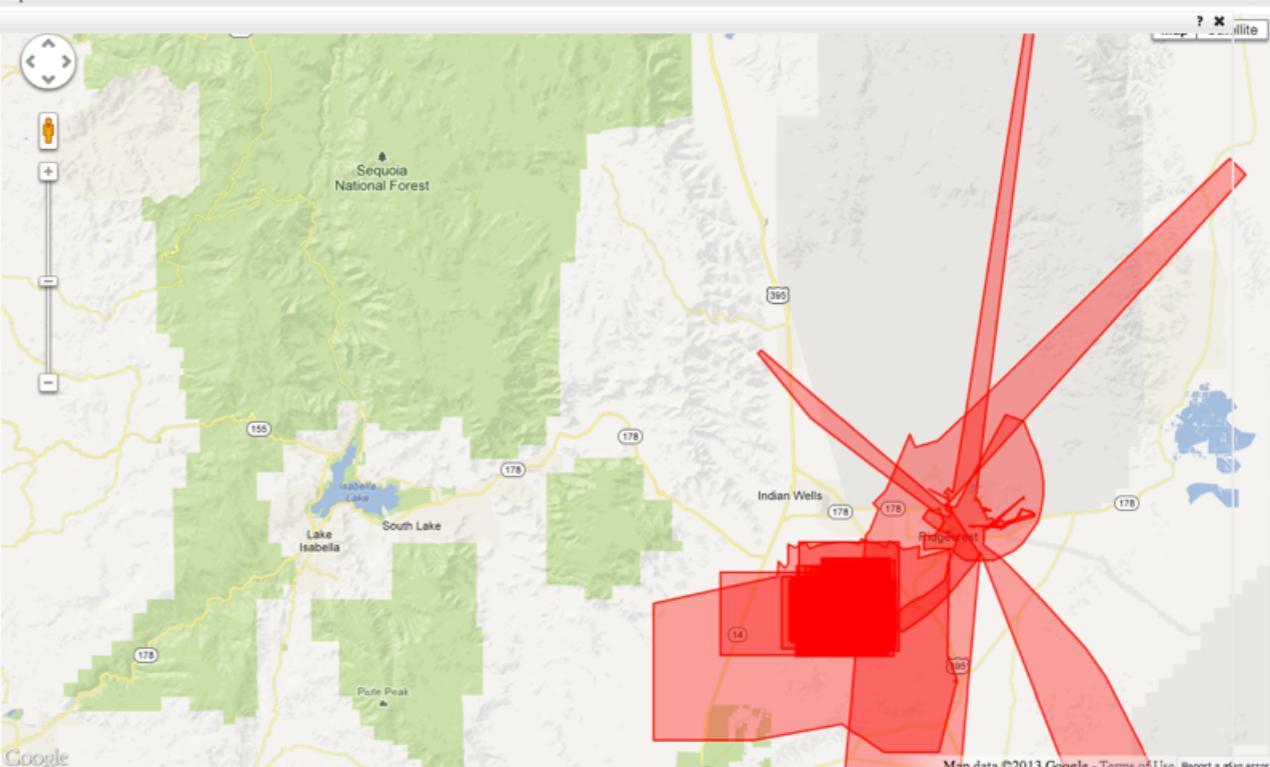
- LIST OF RESULTS
- WEB OF INFORMATION
- IMAGE THUMBNAILS

- MAP DISPLAY
- COLLECTION FOOTPRINTS
- TIMELINE DISPLAY

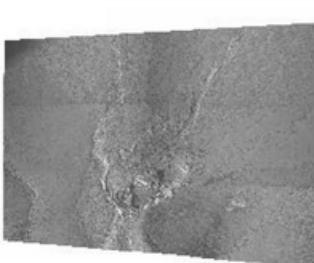
Associations

```
Terms
geometry
Concept : Collection Footprint
Perceivers : FootprintCreator
has Collection Footprint : clip_20120306_193042_967_SEDO.jpg
geometry
Concept : Collection Footprint
Perceivers : FootprintCreator
has Collection Footprint : udpcapture-1330976294561_IR.jpg
geometry
Concept : Image
Perceivers : NITFDigestBean
assoc Image Category : IR
assoc Image Date and Time : 20120229014003
assoc Image Geographic Location : +35.600-117.831+35.600-117.744+35.523-117.744+35.523-117.831
assoc Image Identifier 1 : AF_IR_JPG
assoc Image Identifier 2 : Default Image Identifier
assoc Image Source : AngelFire IR Sensor
geometry
Concept : Collection Footprint
Perceivers : FootprintCreator
has Collection Footprint : udpcapture-1331062280768_HD.jpg
geometry
Concept : Collection Footprint
Perceivers : FootprintCreator
has Collection Footprint : Collection Footprint : clip_20120229_204037_880_UsableData_SCAN_EAGLE.jpg
geometry
Concept : Collection Footprint
Perceivers : FootprintCreator
has Collection Footprint : udpcapture-1330978103845_SDEQ.jpg
geometry
Concept : Collection Footprint
Perceivers : FootprintCreator
has Collection Footprint : udpcapture-1331149334173_BlueDevil HD.jpg
geometry
Concept : Collection Footprint
Perceivers : FootprintCreator
has Collection Footprint : udpcapture-1331149334173_BlueDevil HD.jpg
geometry
Concept : Collection Footprint
Perceivers : FootprintCreator
```

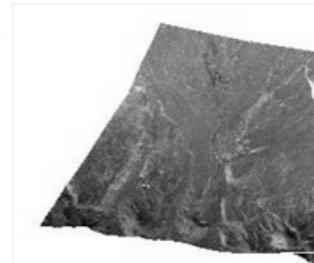
■ **Map**



Imagery



20120301002848-01000834-VIS.ntf
Concept :Thumbnail



20120229013834-05010726-IR.ntf
Concept :Thumbnail



20120229013555-05010608-IR.ntf
Concept :Thumbnail



20120301002848-01000834-VIS.ntf
Concept :Thumbnail

Timeline

Feb 29	<ul style="list-style-type: none"> u0f9bf84-a1c5-4af7-a115-b3ab63ab6d3 2fc1e037-dea1-45bb-83f1-24ca72841dfb u5b53444-7720-aec4-8141-6d7fa4d8cc48 e6c1eae4-47e5-4999-e951-be1fe88ec38 e1698cb8-d664-4026-b571-e2328842d723 67e07078-88cf-4bb8-8f3-618044685428 fbbaedf2-9b99-45b8-9d85-0265deec30e9 2a42d446-4d44-47e0-85b7-d162af12246 3e2c2eb-3274-43b4-be5d-2d1ad28822b8 1e943290-a074-447e-ae54-95db38444c39 29308133-c5ec-46b5-a1e7-062b6ea39b00 56b83b2d-1ac9-455b-9d48-6602fd3dc0d4 d51e03aa-a79e-44ef-8fe-533a03b8564 0d491cc8-941d-4231-b224-e314ac320977 2f60d6be-8d1e-4d05-8e33-b6b63ac59ee8 6e05b067-4590-48b9-a9a5-4344d37b724c a1e1ead2-2002-402d-a8fe-3185cdfff2d 7690ca80-1226-4b6b-akde-d5a0ee877eeb 8740ae9-5d5d-4bdd-88e2-acde8c8a595d 9a3e07e6-70e5-41e5-bb46-e6e4d76c597d bc250562-197a-4de2-8a01-cf0fd1d4ff58 9dc5bfb9-47e9-4800-b87b-9b725c39062 bbe4d57-53c4-45f1-9b99-3dd3e2c0094 190a872d-7358-4546-b523-5e884f111111 	Mar 1
	Feb 29	Mar 1

2012

MISSION FOCUS UNCLASSIFIED // FOUO // COPYRIGHT © 2013 INSTITUTE FOR MODERN INTELLIGENCE, 501(c)3. ALL RIGHTS RESERVED

92

Thursday, August 29, 13

SHAPEFILE FEATURES

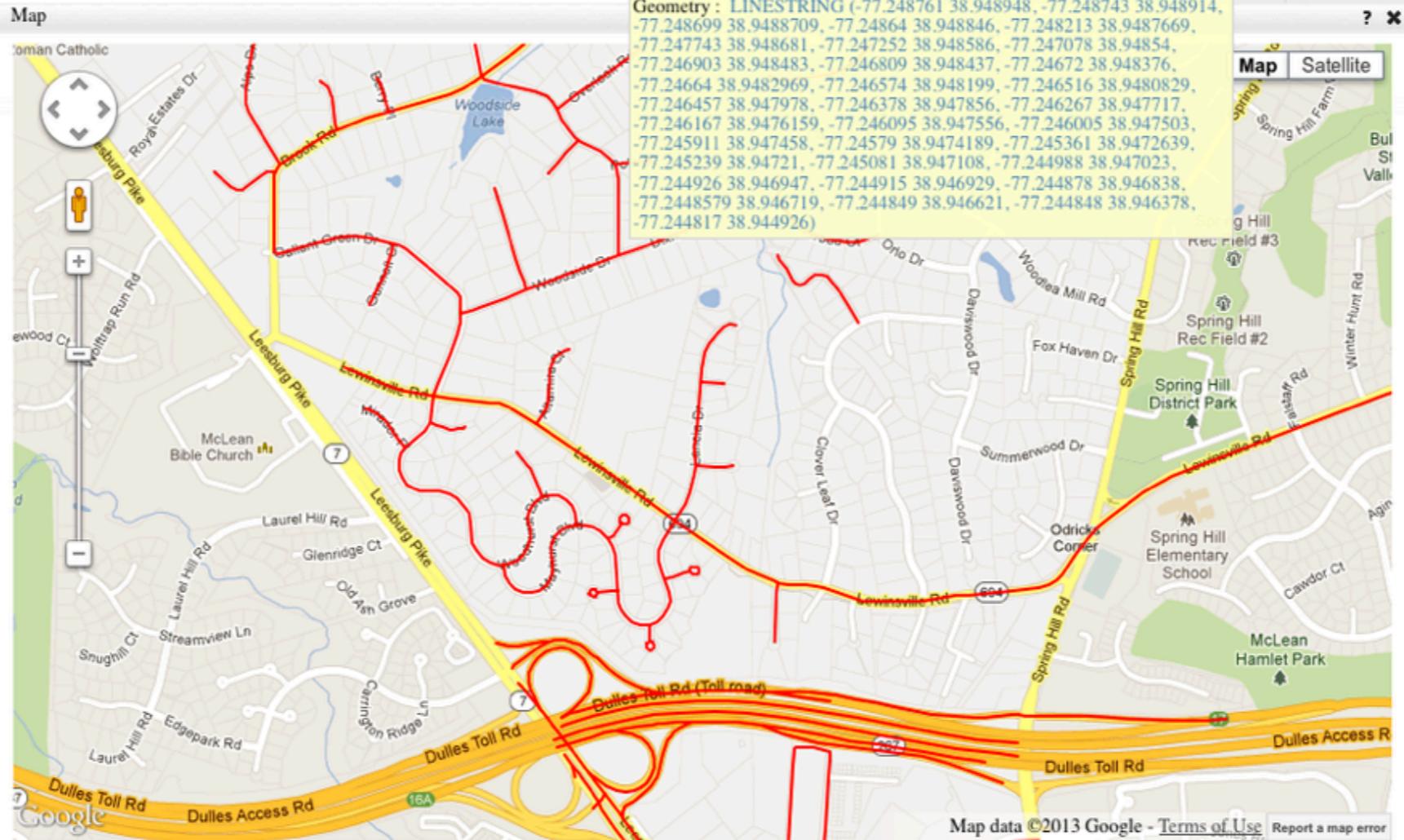
Map | Timeline | Map | Reset

?

Timeline

Sign : Woodside Dr
 Concept : Road Name
 Perceivers : ShapeDigest.district_of_columbia_highway.decoder
 assoc Number of Lanes : *****
 assoc ONEWAY :
 assoc Road Type : residential
 Geometry : LINESTRING (-77.248761 38.948948, -77.248743 38.948914, -77.248699 38.9488709, -77.24864 38.948846, -77.248213 38.9487669, -77.247743 38.948681, -77.247252 38.948586, -77.247078 38.94854, -77.246903 38.948483, -77.246809 38.948437, -77.24672 38.948376, -77.24664 38.9482969, -77.246574 38.948199, -77.246516 38.9480829, -77.246457 38.947978, -77.246378 38.947856, -77.246267 38.947717, -77.246167 38.9476159, -77.246095 38.947556, -77.246005 38.947503, -77.245911 38.947458, -77.24579 38.9474189, -77.245361 38.9472639, -77.245239 38.94721, -77.245081 38.947108, -77.244988 38.947023, -77.244926 38.946947, -77.244915 38.946929, -77.244878 38.946838, -77.2448579 38.946719, -77.244849 38.946621, -77.244848 38.946378, -77.244817 38.944926)

Map | Satellite



Get Current Map Results as KML
 View Current Map Results in Google Earth
 Get Current Map Results as KMZ

Map data ©2013 Google - Terms of Use Report a map error

Get All Search Results as KML
 View All Search Results in Google Earth
 Get All Search Results as KMZ

CAPABILITIES - OTHER

■ MONITORING & METERING

- PERSIST LOG MESSAGES IN THE SIGN-SPACE
- MONITOR & METER
 - DISK & CPU UTILIZATION ACROSS THE CLUSTER
 - HADOOP & ACCUMULO SW INFRASTRUCTURE
 - SEMAPP CORE METHODS AND SERVICES
 - INGEST MESSAGE TRANSFERS

■ SECURITY

- DATA CLASSIFICATION PERSISTED AT THE CELL-LEVEL
- RELEASABILITY DETERMINED BY USER CREDENTIALS & DATA CLASSIFICATION
- APPLICATIONS USER AUTHENTICATION & AUTHORIZATION
- ENGINEERED FOR PL4

■ CONTINUOUS INTEGRATION

- UNIT TESTS (AS RESULTS FROM TEST DRIVEN DEVELOPMENT)
- TEST FRAMEWORK - SET SIGN-SPACE STATE, MANAGE TEST SUITES
- MAVEN BUILD & DEPENDENCY MANAGEMENT
- LEININGEN & NODEJS WEB TEST & DEPENDENCY MANAGEMENT
- GIT CODE & ARTIFACT REPOSITORIES
- JENKINS AUTOMATED BUILD, QUALITY CHECKS, & REGRESSION TEST
- AMBIENT VISUALIZATION OF BUILD & TEST STATUS

■ DEPLOYMENT - INSTALL, CONFIG, TEST ON PRODUCTION HARDWARE

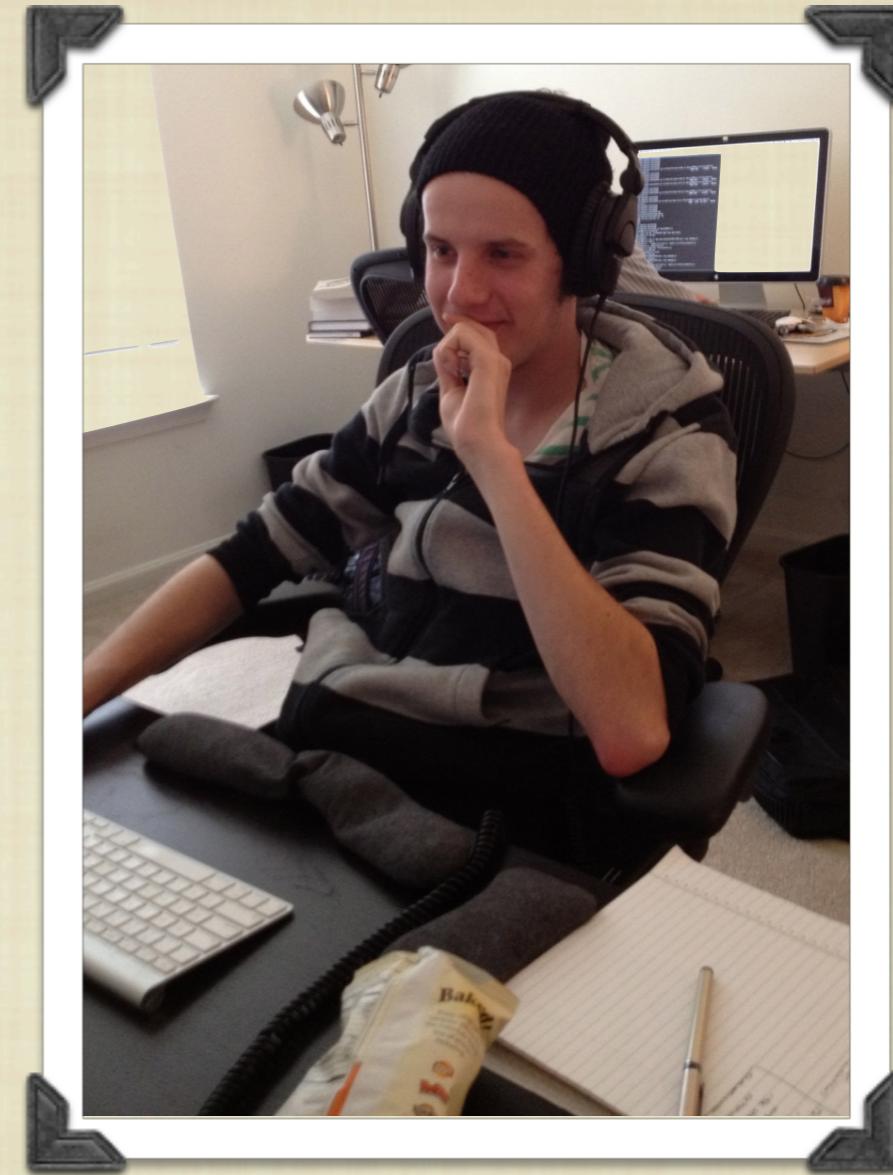
- SEMAPP KICKSTART
 - RHEL + PACKAGES INSTALL (STIG'D)
 - CENTOS
 - DISK PARTITIONING & NIC BONDING
- PROCEDURES FOR NTP CONFIG, CLOUD INFRASTRUCTURE AND SEMAPP SOFTWARE DEPLOY
- TEST PROCEDURES FOR CONFIGURATION AND CORE FUNCTIONALITY
- PASSED AREA 59 FAT



TEST DATA

COLLECTED FROM THESE SOURCES

- ALOE 1 & 2
 - IMAGERY, VIDEO, WAMI, CHAT
- BLUE DEVIL ONLINE ARCHIVE
 - IMAGERY, VIDEO
- MICHELSON LAB ARCHIVES, CHINA LAKE
 - VARIOUS
- OPEN SOURCE
 - GOVTRACK
 - TWITTER
 - WIKIPEDIA
 - FEDERAL CONTRACT SPENDING
 - AUDIO METADATA
 - US PATENTS
 - PROJECT GUTENBERG
 - OPEN SECRETS
 - OPENSTREETMAP
- SAIC GSTI
 - LIDAR
 - HSI



DEVELOPMENT BACKLOG



Move it forward. Make it happen.



ONGOING WORK

- **ASSERT, RETRIEVE, DELETE, EXISTS, TALLY**
 - GEOMETRY EXTENSIONS
 - SIGNS WITH MULTIPLE GEOMETRIES
- **API FOR VIDEO PLAYBACK**
- **APPLICATIONS**
 - DATA MANAGEMENT
 - FOUNDATION DATA BUILDER
 - EXECUTABLE PROCESS LOADER
 - MI PLAYER
 - OPERATORS CONSOLE
 - SEMAPP NOTEBOOK
 - SEMAPP SUITE
 - TRACKING
 - USER ASSERTS
 - USER SA
 - USER SELF-REGISTRATION
 - ...
- **CAPACITY & STABILITY**
 - CACHING
 - CUCUMBER
 - LOAD BALANCING
- **DEPLOYMENT ACTIVITIES**
 - DEPLOYMENT READINESS
 - IA
 - SITE INTEGRATION
 - THEATER SUPPORT
- **DOCUMENTATION**
 - APPLICATION USER GUIDES
 - CONOPS
 - DATA ARCHITECTURE GUIDE
 - DEVELOPER'S GUIDE
 - MANAGEMENT GUIDE
 - OPERATOR'S GUIDE
 - PROCESS DOCUMENTATION
 - TRAINING MATERIALS
- **INGEST**
 - DIRECT FROM RDB
- **INTEGRATION**
 - BULK IMPORT / EXPORT
 - DIB
 - JPIP
 - NVS-ABI
 - SOM
 - VIRGO
 - WORLDWIND
- **M&M QUERY PERFORMANCE**
- **MODEL SERVICES**
 - QUERY BROADEN / NARROW
 - SAMEAS HANDLING
- **PROCESS**
 - ADVANCED REASONING ENGINES
 - IMAGE PROCESSING
 - MESSAGE CHAINING
 - MICROSOFT PRODUCTS
 - NLP ANALYTICS
 - RDF
 - USER-DEFINED M/R
 - ...
- **MIXED MODALITY COMPOUND QUERIES**

Move it forward. Make it happen.

RECENT RELEASE RECORD



Move it forward. Make it happen.

RELEASE 0.13 FOCUS

How do we know it's working?

How do we keep it working as it grows?

TESTING

PERFORMANCE

QUALITY

■ TEST CASE DEVELOPMENT

- FUNCTIONAL
 - DATA & INDEX SERVICES
 - PROCESS SERVICES
 - PROCESSES
 - QUERY

■ PERFORMANCE

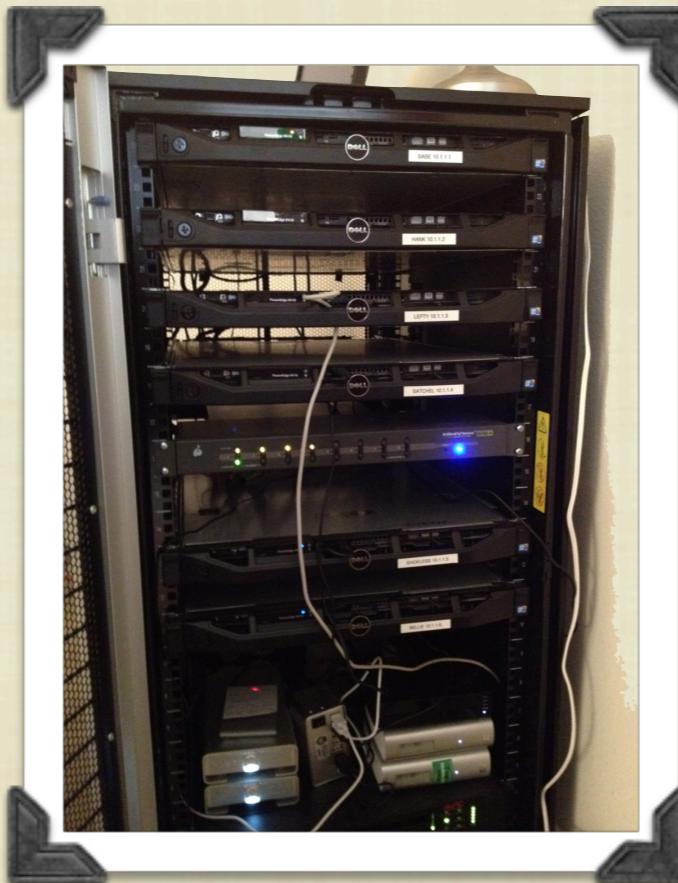
- DATA & INDEX SERVICES
- INGEST
- PROCESSING
- QUERY

■ CONFIGURATION

- HADOOP, ZOOKEEPER, ACCUMULO
- SEMAPP CORE
- M&M
- ENTERPRISE

■ ACCEPTANCE

- SEARCH - SURF APPLICATION
- COMMAND-LINE APPLICATIONS
- WEB SERVICES



■ TEST MANAGEMENT

- ORGANIZATION - TEST CASES, SETS, & PLANS
- DOCUMENTATION & CM
- CONTINUOUS INTEGRATION
- MANUAL TESTING PROCEDURES

■ CODE QUALITY

- CODING STANDARDS COMPLIANT THROUGHOUT
- CONSISTENT USE OF PATTERNS
- SCALABLE
- THREAD SAFE
- USEFUL LOG / ERROR MESSAGES

■ MONITORING & METERING

- OPERATOR MONITORING CONSOLE
- CORE HEALTH & STATUS INDICATORS
- AGGREGATED PERFORMANCE METRICS

■ PRACTICE

- DEPLOY IT
- STRESS IT
- TREAT IT LIKE A PRODUCTION SYSTEM

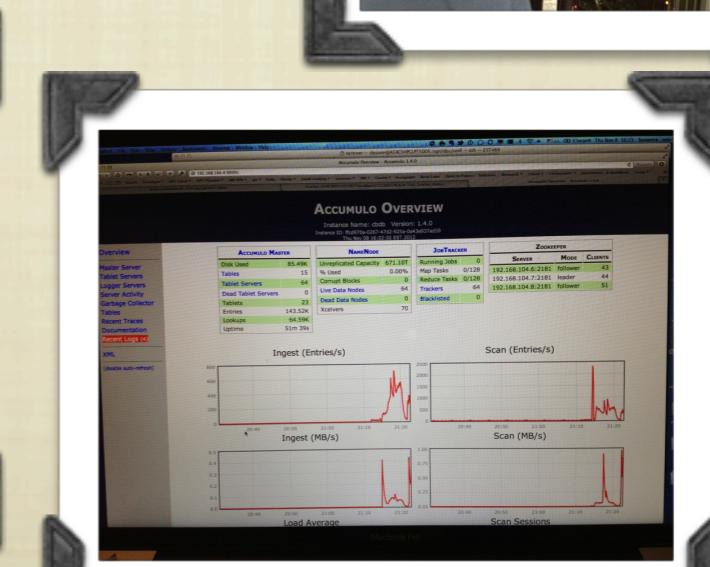
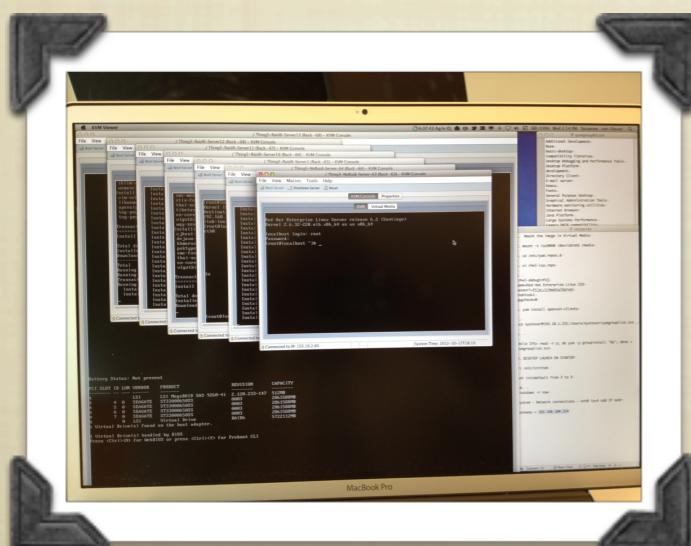
RELEASE 0.14 FOCUS

DEPLOY ONTO PRODUCTION HARDWARE

The “Paradigm Cluster”

ON ALL 64 DATA NODES AND 15 SERVERS

- FIX INITIAL RHEL INSTALL
- DEVELOP SEMAPP KICKSTART
- PARTITION THE DRIVES
- BOND THE NICs
- CONFIGURE NTP
- RESOLVE HARDWARE FAILURES
- WORK THROUGH UCSM ISSUES
- IA SCAN
- INSTALL, CONFIGURE, AND TEST CLOUD INFRASTRUCTURE
- INSTALL, CONFIGURE, AND TEST SEMAPP SOFTWARE
- DOCUMENT / UPDATE DEPLOYMENT PROCEDURES
- ORGANIZE AND REFINE TEST PROCEDURES FOR CONFIGURATION AND CORE FUNCTIONALITY
- PARTICIPATE IN AREA 59 FAT





RELEASE 0.15 FOCUS

Administrator, Operator, User Support

- THEATER CONOPS
- THEATER WORK FORCE
- NVS INTEGRATION
- APPLICATIONS
 - SEARCH-SURF
 - DATA MANAGEMENT
 - INGEST
 - PROCESS MANAGEMENT
- DOCUMENTATION
 - APPLICATIONS
 - OPERATION
 - MANAGEMENT
 - SYSTEM SERVICES / APIs
 - TRAINING

—USABILITY—

Move it forward. Make it happen.



RELEASE 0.16 FOCUS

— DOMAIN INTEGRATION —

- ACTIVITY BASED INTELLIGENCE
- BROADEN SCOPE OF DATA
- A FRAMEWORK FOR WEB APPLICATIONS

Getting together with domain SMEs

Move it forward. Make it happen.



RELEASE 0.17 FOCUS

— FRONT END —

- LOOK & FEEL
- APPLICATION PROTOTYPE(S)
- INFORMATION VISUALIZATION
- MORE GEOINT DATA TYPES

Move it forward. Make it happen.



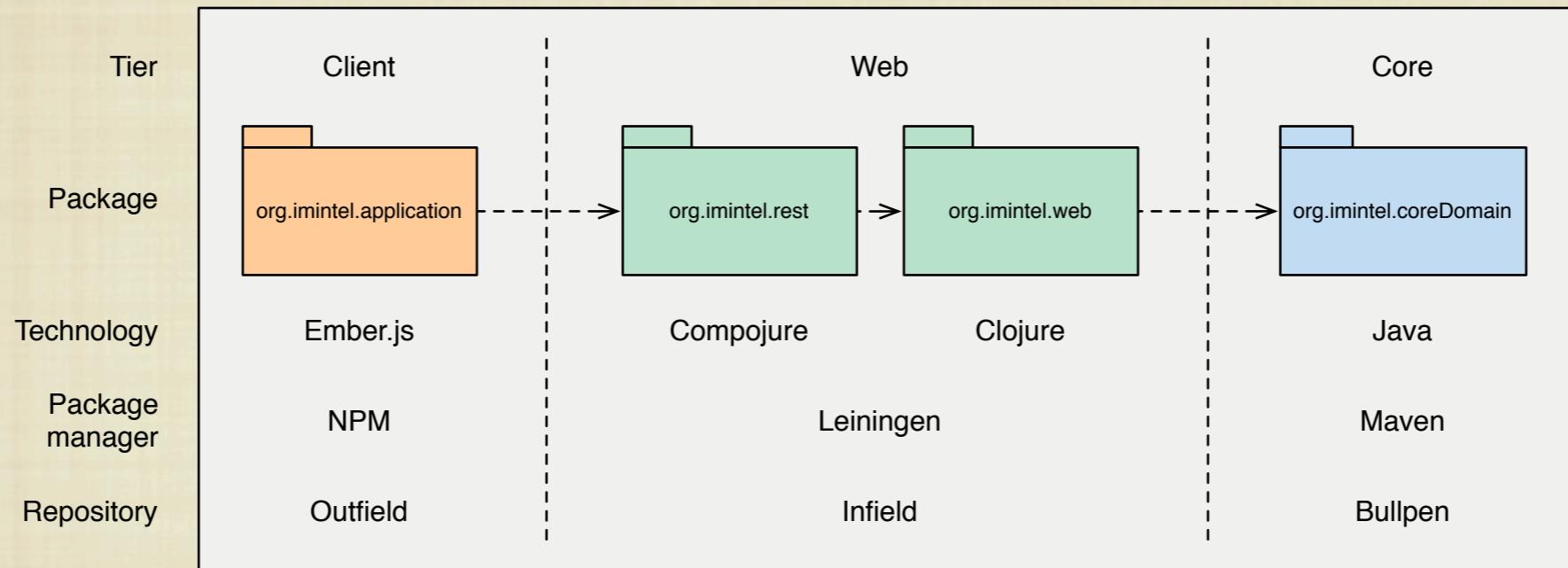
RELEASE 0.18 FOCUS

NEW APPS

■ THE EYE

■ DATA SA

Move it forward. Make it happen.

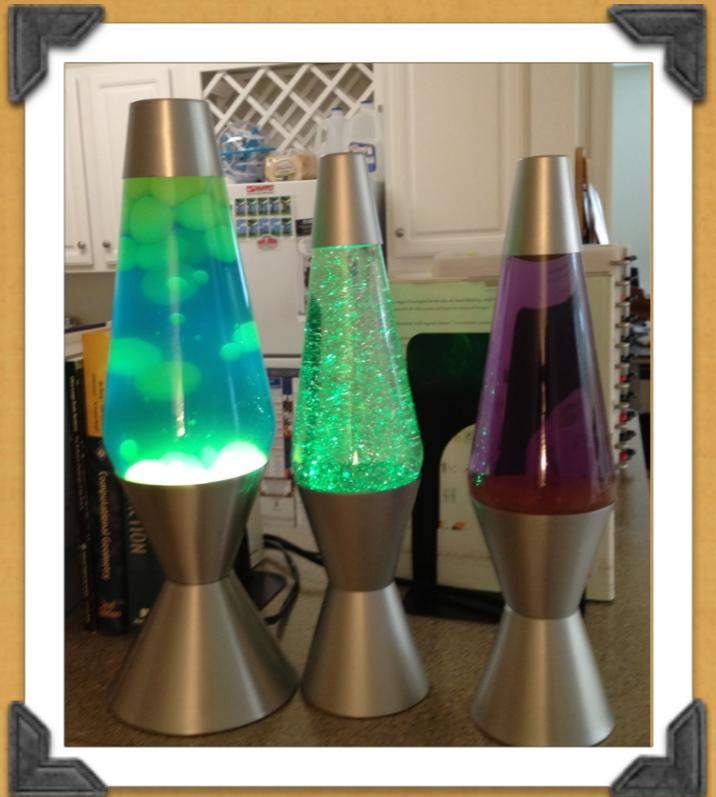




LOOKING FORWARD

→ RELEASE IT! ←

- DEPLOY TO LABORATORY / OPERATIONAL (LABOP) ENVIRONMENT(S)
- PURSUE CERTIFICATION AND ACCREDITATION TO OPERATE
- INTEGRATE WITH LABOP SERVICES (E.G. MONITORING)
- INGEST LABOP DATA
- PROCESS DATA
- WORK WITH USERS (E.G. ANALYSTS) AND LABOP OPERATORS
- INTEROPERATE WITH ENTERPRISE SERVICES (E.G. GEOBUS, DIB)
- REFINE / EXTEND / DEVELOP APPLICATIONS, INTERFACES, SERVICES, PROCESSES, AND DOMAIN-MODELS AS NEEDED TO SUPPORT LABOP OBJECTIVES
- SUPPORT AND MAINTAIN SYSTEM OPERATION



RELEASE PLANS



Move it forward. Make it happen.

R0.13 - R0.19 Schedule

— Obsolete —

	Sept-Oct 2012	Nov-Dec 2012	Jan-Feb 2013	Mar-Apr 2013	May-Jun 2013	Jul-Aug 2013	Sept-Oct 2013
	R 0.13	R 0.14	R 0.15	R 0.16	R 0.17	R 0.18	R 0.19
Core Services & Processes	Gazetteer JPIP Shape File M&M Performance	Virgo Integration NVS Integration SensrWeb Integration	Query Broaden Query Narrow RDB Ingest Named Entity	SameAs Handling Walk Back Model Export Query Bool Ext	Bulk Export SDA Improvements	Search Engine NVS ABI Integration	User Messaging Automated Reasoning
Applications	Unif Geo Extr File Std Det Monitor	System Mgmt	Process Reg User Self-Reg	User Assert User Ingest	SemApp Notebook Data SA	Search Engine ABI	User SA Mission Scrum
Test & IA	Func, Perf, Conf Tests Test Plans Test Mgmt Acceptance Tests	IV&V	Theater Smoke Test Regression Testing	Regression Testing Func, Perf, Conf Tests Test Plans Acceptance Tests			
Deployment & O&M	ConUS Install&Config	IA Scan&Remediate	Theater Deploy Site Integration	Theater Support Config & SW Update New Feature Deploy			
Training	Process Docs	Operator's Guide Priv Usr App Guides	Process Reg Guide User Self-Reg Guide Operator Training	User Assert Guide User Ingest Guide User Training	Developers Guide Data Arch Guide Developer Training	Operator Training User Training Developer Training	Operator Training User Training Developer Training
Completed In progress Defined							

R0.15 - R0.21 Schedule

— Obsolete —

	Jan-Feb 2013	Mar-Apr 2013	May-Jun 2013	Jul-Aug 2013	Sept-Oct 2013	Nov-Dec 2013	Jan-Feb 2014
Core Services & Processes	R 0.15	R 0.16	R 0.17	R 0.18	R 0.19	R 0.20	R 0.21
Applications	Gazetteer NVS Integration Shape File Index decoupling	ABI Integration Virgo Integration Index associations OWL expressiveness	Bulk Export Model Export RDB Ingest	SameAs Handling Query Bool Ext	Query Broaden Query Narrow Search Engine	Named Entity JPIP	User Messaging Automated Reasoning
Test & IA	Web app security Data Mgmt Process Mgmt	Operator Console	User Assert User Ingest	SemApp Notebook Data SA	ABI	Search Engine User SA	Process Reg Mission Scrum
Deployment & O&M	Performance Testing	Interface Testing	Theater Smoke Test Regression Testing IA Scan	Regression Testing Theater ATO	Regression Testing	Regression Testing	Regression Testing
Documentation & Training	Deployment prep	Deployment prep	China Lake Update	Theater Deploy Site Integration	Theater Support Config & SW Update New Feature Deploy	Theater Support Config & SW Update New Feature Deploy	Theater Support Config & SW Update New Feature Deploy
	Operator's Guide Management Guide Process Docs REST API Doc	ConOps	Data Arch Guide Priv User App Guides Search-Surf Guide	Operator Training Manager Training User Assert Guide User Ingest Guide	User Training Developers Guide Notebook Guide Data SA Guide	Developer Training Training (as needed) ABI Guide	Training (as needed) User SA Guide Mission Scrum Guide
Completed In progress Slow / Blocked							

Ro.15 - Ro.21 Schedule

	May-Jun 2013	Jul-Aug 2013	Sept-Oct 2013	Nov-Dec 2013	Jan-Feb 2014	Mar-Apr 2014
Core Services & Processes	R 0.17	R 0.18 Now	R 0.19	R 0.21	R 0.22	R 0.23
Applications	Data SA					
Other	Office move Dev ref refactor					
Completed						
In progress						
Slow / Blocked						

Off Contract

Scale changes everything

Walk with your head in the clouds
and your feet in the dirt

Everything changes when you
deploy to a production
environment

The ground state of
data is not a pure
crystalline form

Everything changes when
you start writing the code

There's More

In_all_of_it_Together
Than we are getting or
could possibly get out of
All of the Parts

Don't sell yourself short (the rest of the world will do it for you soon enough)

Our universe of data & processing is an
Ultra-Large Scale system

New science waits to be discovered

Simple is hard

Diversity gives Intel its richness

Every system should
have the ability to
export its data

Our data & processing assets
belong to the nation
(not the system or the contractor)

What you don't have
in your head, you gotta
have in your feet

The sink is not a
repository for dirty dishes

Godspeed

Be happy, write code!

