# UC San Diego
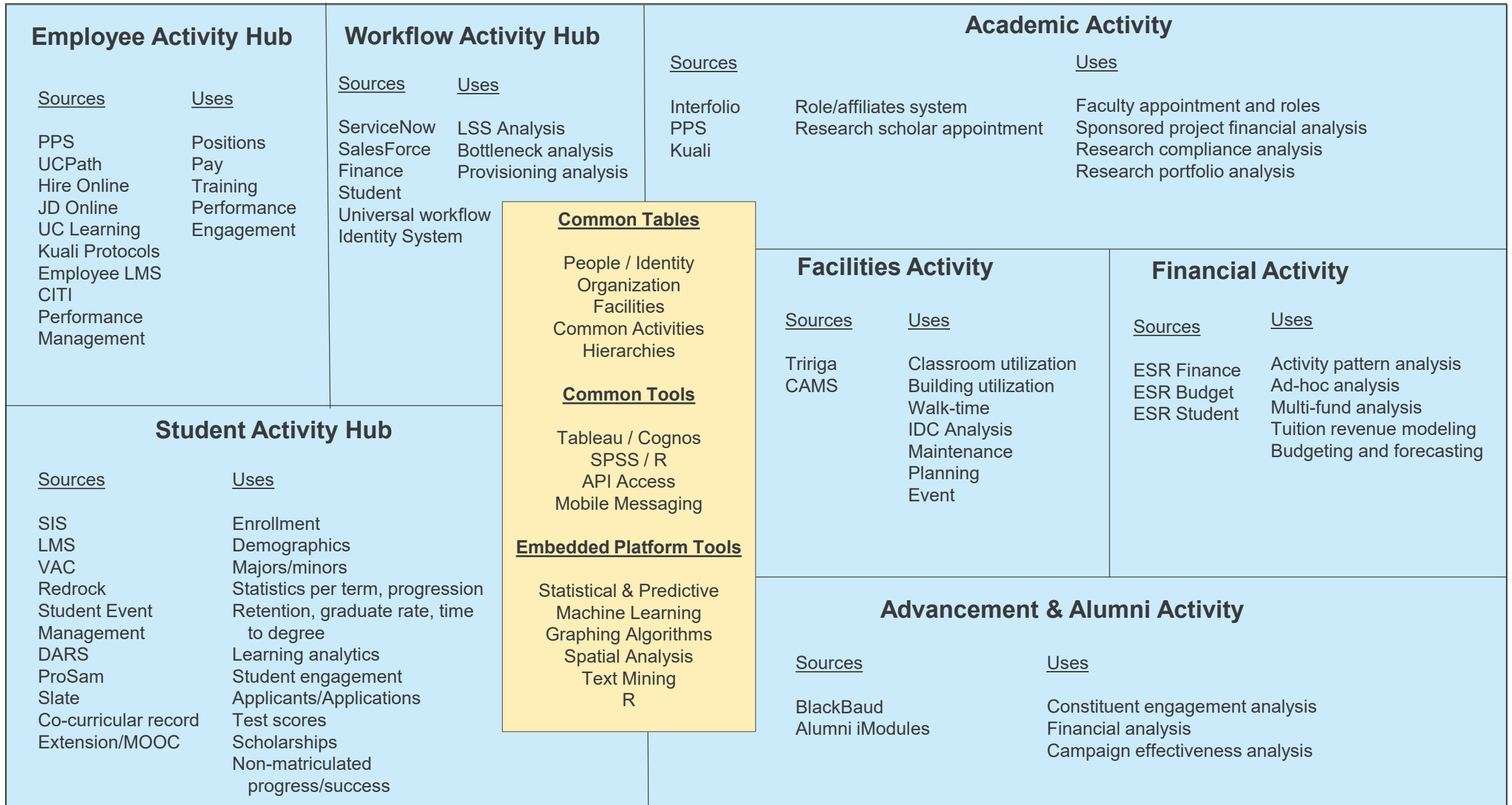
## ACTIVITY HUBS

Next generation data and analytics platform for the campus

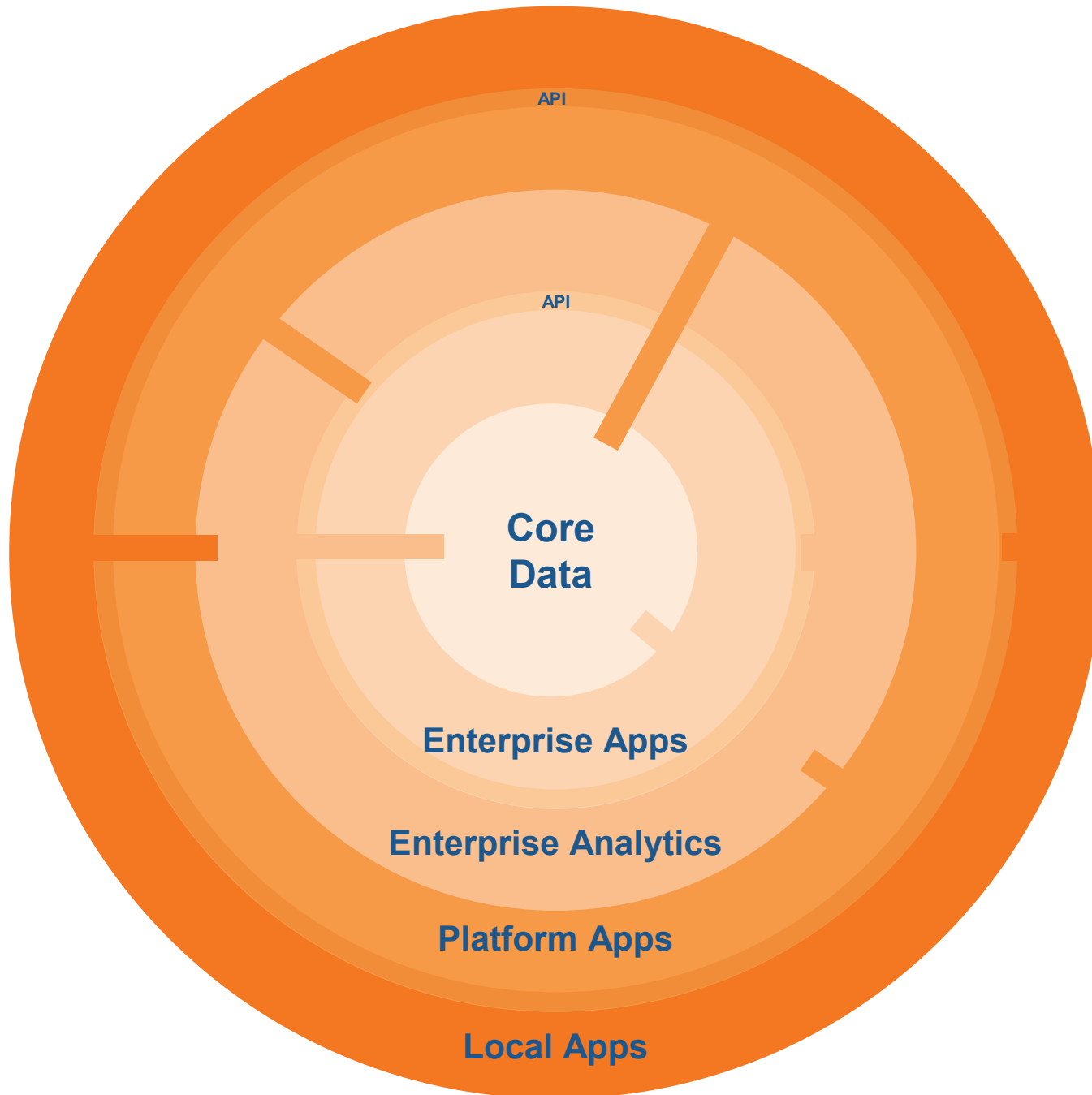June 20, 2018

*This is a living document subject to substantial revision!*

# Overview of the next generation data warehouse

## Employee Activity Hub

**Sources**

PPS
UCPath
Hire Online
JD Online
UC Learning
Kuali Protocols
Employee LMS
CITI
Performance
Management

**Uses**

Positions
Pay
Training
Performance
Engagement

## Workflow Activity Hub

**Sources**

ServiceNow
SalesForce
Finance
Student
Universal workflow
Identity System

**Uses**

LSS Analysis
Bottleneck analysis
Provisioning analysis

## Academic Activity

**Sources**

Interfolio
PPS
Kuali

Role/affiliates system
Research scholar appointment

**Uses**

Faculty appointment and roles
Sponsored project financial analysis
Research compliance analysis
Research portfolio analysis

## Common Tables

People / Identity
Organization
Facilities
Common Activities
Hierarchies

## Common Tools

Tableau / Cognos
SPSS / R
API Access
Mobile Messaging

## Embedded Platform Tools

Statistical & Predictive
Machine Learning
Graphing Algorithms
Spatial Analysis
Text Mining
R

## Facilities Activity

**Sources**

Tririga
CAMS

**Uses**

Classroom utilization
Building utilization
Walk-time
IDC Analysis
Maintenance
Planning
Event

## Financial Activity

**Sources**

ESR Finance
ESR Budget
ESR Student

**Uses**

Activity pattern analysis
Ad-hoc analysis
Multi-fund analysis
Tuition revenue modeling
Budgeting and forecasting

## Student Activity Hub

**Sources**

SIS
LMS
VAC
Redrock
Student Event
Management
DARS
ProSam
Slate
Co-curricular record
Extension/MOOC

**Uses**

Enrollment
Demographics
Majors/minors
Statistics per term, progression
Retention, graduate rate, time
   to degree
Learning analytics
Student engagement
Applicants/Applications
Test scores
Scholarships
Non-matriculated
   progress/success

## Advancement & Alumni Activity

**Sources**

BlackBaud
Alumni iModules

**Uses**

Constituent engagement analysis
Financial analysis
Campaign effectiveness analysis

# ESR 'layered' architecture

**Local applications**
- Support local innovations and needs, can come and go
- Access data, authenticate via APIs
- Are directly access by end users
- Can be promoted to the platform or enterprise application layer
- Need to respect scale and security standards

**Platform applications**
- Have a well-designed real-time API architecture
- APIs are used by central/distributed IT staff from local apps
- Provide access, manage workflows, content, collaboration
- Span multiple business functions
- Endure and evolve
- Enterprise content management, workflow, IAM tools, visualization

**Enterprise analytics**
- Are independent of enterprise applications
- Access data directly, can consume APIs
- Are directly accessed by end users
- Can be embedded in local apps
- Institutional data warehouses, data lakes
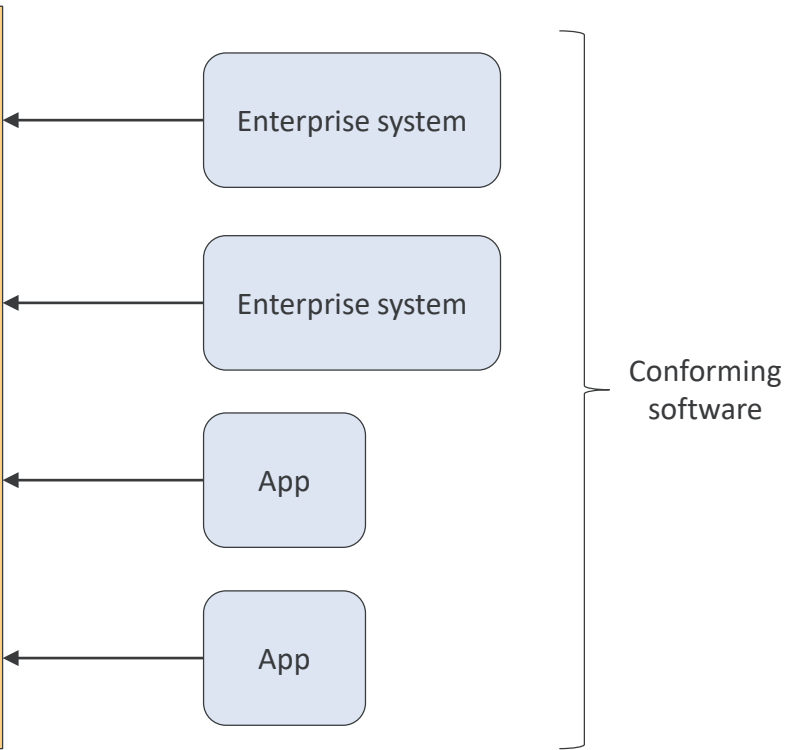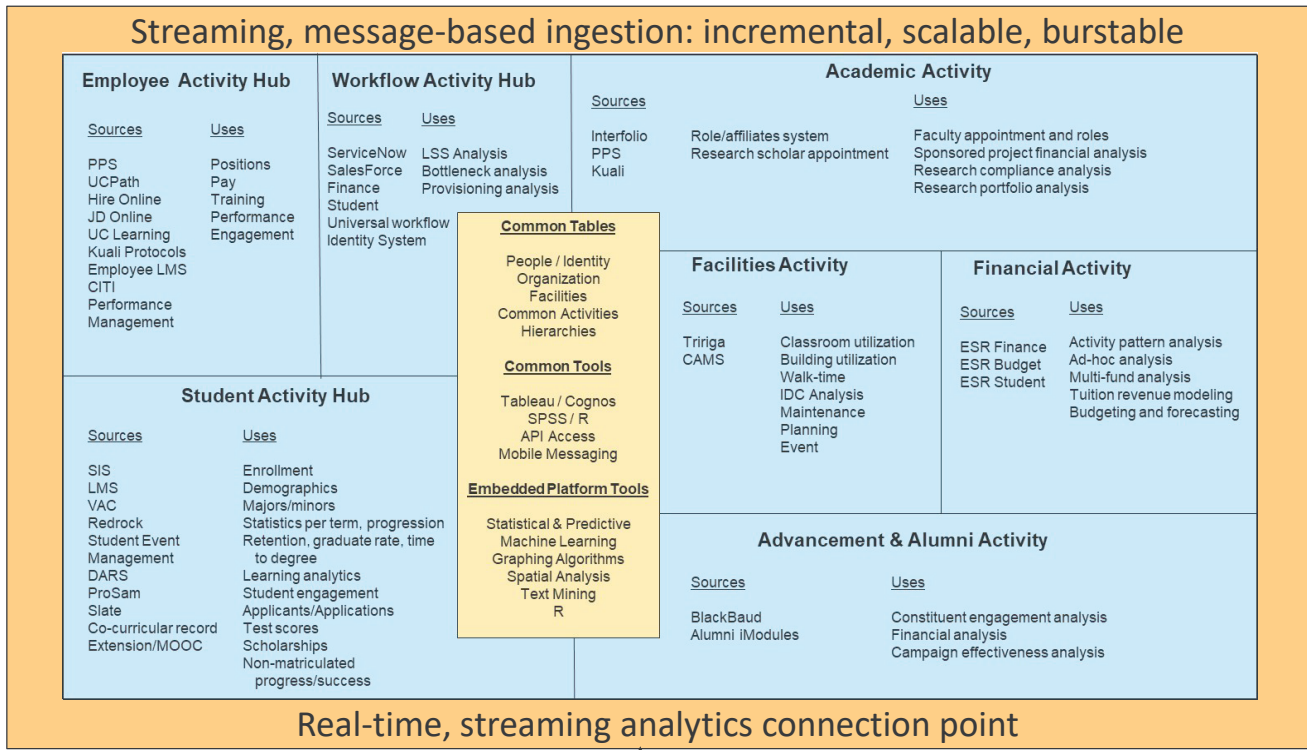
**Enterprise applications**
- Have a real-time transaction API layer available to upper layers
- Endure, are rarely replaced
- Manage core transactions, align with specific business functions
- Are directly access by end users
- May have analytics within and can access core data
- Finance, student information, HR, budgeting systems

**Core Data**
- Data essential to all layers
- Core transaction and master data
- Does not include local application data

3

# Activity hubs and new enterprise systems



4

## Streaming, message-based ingestion: incremental, scalable, burstable

**Employee Activity Hub**

Sources
- PPS
- UCPath
- Hire Online
- JD Online
- UC Learning
- Kuali Protocols
- Employee LMS
- CITI
- Performance Management

Uses
- Positions
- Pay
- Training
- Performance
- Engagement

**Workflow Activity Hub**

Sources
- ServiceNow
- SalesForce
- Finance
- Student
- Universal workflow
- Identity System

Uses
- LSS Analysis
- Bottleneck analysis
- Provisioning analysis

**Academic Activity**

Sources
- Interfolio
- PPS
- Kuali

Uses
- Role/affiliates system
- Research scholar appointment
- Faculty appointment and roles
- Sponsored project financial analysis
- Research compliance analysis
- Research portfolio analysis

**Common Tables**

- People / Identity
- Organization
- Facilities
- Common Activities
- Hierarchies

**Common Tools**

- Tableau / Cognos
- SPSS / R
- API Access
- Mobile Messaging

**Embedded Platform Tools**

- Statistical & Predictive
- Machine Learning
- Graphing Algorithms
- Spatial Analysis
- Text Mining
- R

**Facilities Activity**

Sources
- Tririga
- CAMS

Uses
- Classroom utilization
- Building utilization
- Walk-time
- IDC Analysis
- Maintenance
- Planning
- Event

**Financial Activity**

Sources
- ESR Finance
- ESR Budget
- ESR Student

Uses
- Activity pattern analysis
- Ad-hoc analysis
- Multi-fund analysis
- Tuition revenue modeling
- Budgeting and forecasting

**Student Activity Hub**

Sources
- SIS
- LMS
- VAC
- Redrock
- Student Event Management
- DARS
- ProSam
- Slate
- Co-curricular record
- Extension/MOOC

Uses
- Enrollment
- Demographics
- Majors/minors
- Statistics per term, progression
- Retention, graduate rate, time to degree
- Learning analytics
- Student engagement
- Applicants/Applications
- Test scores
- Scholarships
- Non-matriculated progress/success

**Advancement & Alumni Activity**

Sources
- BlackBaud
- Alumni iModules

Uses
- Constituent engagement analysis
- Financial analysis
- Campaign effectiveness analysis

## Real-time, streaming analytics connection point

Enterprise system

Enterprise system

App

App

Conforming software

## "Wedge" diff calculator, streaming converter

ODS

ODS

ODS

ODS

Enterprise system

Enterprise system

App

App

Non-conforming software

Activity hubs ingest data via a streaming message service. Curated views and activity tables should employ "duplicate safe" rendering methods, allowing for indempotent messages. This relaxes data consistency significantly, easing the integration complexity.

The streaming analytics connection point allows for directly connecting the streaming ingestion engine with a real-time streaming analytics machine learning platform to process inbound messages
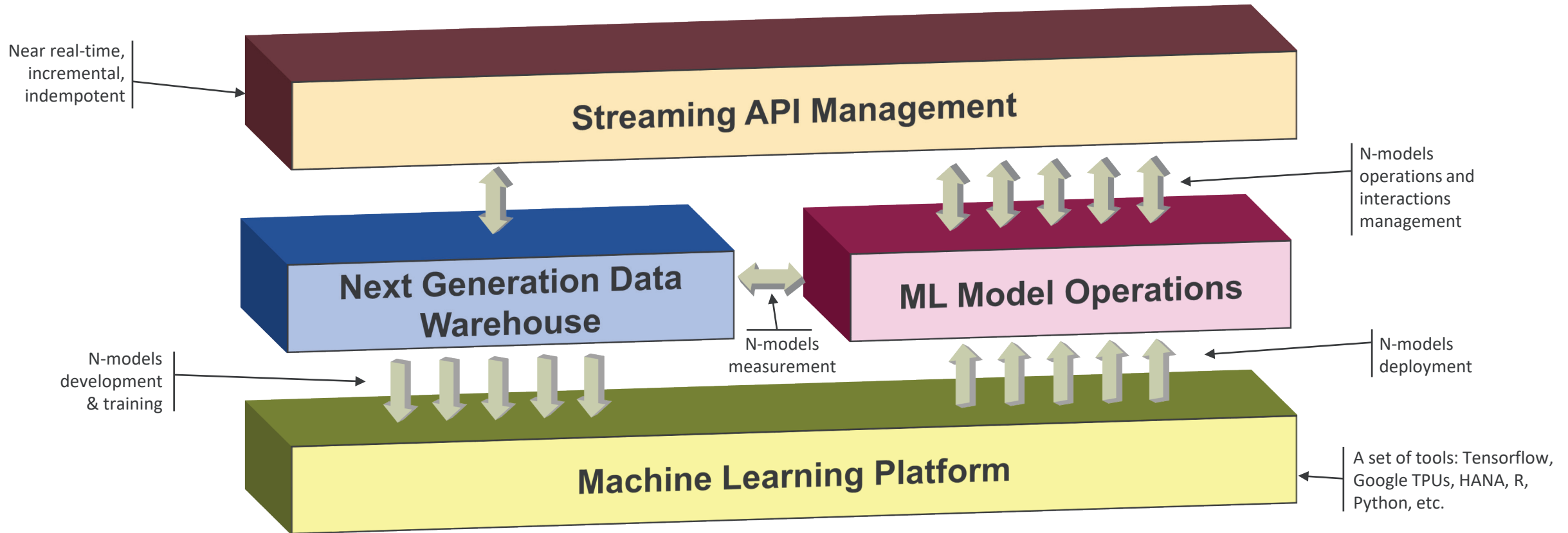
Conforming software meets the streaming message-based ingestion method and submit directly to the activity hub message layer.

Non-conforming software needs a "wedge" integration point that helps calculate differences in snapshots to determine incremental adds, updates and deletes. The ODS and other tools for this wedge can exist in any platform(s), including HANA. The principle define choice is long-term cost and performance needs.

# Managing N-ML models in the next generation analytics environment

How can we use machine learning to improve administrative processes, student success?

- Multiple models may be active per each business opportunity (e.g., student learning feedback, student success intervention, financial activity fraud detection)
- Multiple models will be developed and trained based on prior streams of data
- Multiple models will be deployed to actively interact with real-time streams of data, interacting with requesting systems and users, activating workflows
- Multiple models can be managed within a single pane of glass. Operations can ensure reliability, detect anomalies, bring up and take down models
- Model measurement data feeds back into the next generation data warehouse to guide further model development

Near real-time, incremental, indempotent

**Streaming API Management**

N-models operations and interactions management

**Next Generation Data Warehouse**

N-models measurement

**ML Model Operations**

N-models development & training

N-models deployment

**Machine Learning Platform**

A set of tools: Tensorflow, Google TPUs, HANA, R, Python, etc.

# SAP HANA
## Data Ingestion and Analytics modelling overview

**Consume**

Tableau | Tableau Web Server | Cognos | SPSS | SAS | R | Microsoft Office

**Compute & Data Store**



SAP HANA

Modelling

In-Memory

Views

Procedures

SQL

CDS

Virtual Tables

Flowgraphs

**Data Lake**

SAP Vora

hadoop

Spark

Hortonworks

cloudera

MAPR

altiscale

**Ingest**

ETL ⇧ Replication ⇧ Streaming ⇧ Virtual Access · · ·

**Sources**

SAP S/4HANA     TERADATA     Twitter     Sensor     Machine     IBM DB2     Microsoft SQL Server     ORACLE · · ·

GOOGLE BIGQUERY

# Platform predictive capabilities

## Classification Analysis
- CART
- C4.5 Decision Tree Analysis
- CHAID Decision Tree Analysis
- K Nearest Neighbour
- Logistic Regression Elastic Net
- Back-Propagation (Neural Network)
- Naïve Bayes
- Support Vector Machine
- Random Forests
- Gradient Boosting Decision Tree
- Linear Discriminant Analysis (LDA)
- Confusion Matrix
- Area Under Curve (AUC)
- Parameter Selection/Model Evaluation

## Regression
- Multiple Linear Regression Elastic Net
- Polynomial, Exponential, Bi-Variate Geometric, Bi-Variate Logarithmic Regression
- Generalized Linear Model
- Cox Proportional Hazards Model

## Cluster Analysis
- ABC Classification
- DBSCAN
- K-Means/Accelerated K-Means
- K-Medoid Clustering
- K-Medians
- Kohonen Self-Organized Maps
- Agglomerate Hierarchical
- Affinity Propagation
- Latent Dirichlet Allocation (LDA)
- Gaussian Mixture Model (GMM)
- Cluster Assignment

## Time Series Analysis
- Single/Double/Brown/Triple Exponential Smoothing
- Forecast Smoothing
- Auto – ARIMA/ Seasonal ARIMA
- Croston Method
- Forecast Accuracy Measure
- Linear Regression with Damped Trend and Seasonal Adjustment
- Test for White Noise, Trend, Seasonality
- Fast Fourier Transform (FFT)
- Correlation Function

## Association Analysis
- Apriori
- Apriori Lite
- FP-Growth
- KORD – Top K Rule Discovery
- Sequential Pattern Mining

## Probability Distribution
- Distribution Fit/Weibull analysis
- Cumulative Distribution Function
- Quantile Function
- Kaplan-Meier Survival Analysis

## Outlier Detection
- Inter-Quartile Range Test (Tukey's)
- Variance Test
- Anomaly Detection
- Grubbs Outlier Test

## Recommender
- Factorized Polynomial Regression Models

## Link Prediction
- Common Neighbors
- Jaccard's Coefficient
- Adamic/Adar
- Katz$\beta$

## Statistical Functions
- Mean, Median, Variance, Standard Deviation, Kurtosis, Skewness
- Covariance Matrix
- Pearson Correlations Matrix
- Chi-squared Tests:
    - Test of Quality of Fit
    - Test of Independence
- F-test (variance equal test)
- Data Summary
- ANOVA
- One-sample Median Test
- T Test
- Wilcox Signed Rank Test

## Data Preparation
- Sampling
- Binning
- Scaling
- Partitioning
- Principal Component Analysis (PCA)/ PCA Projection

## Other
- Weighted Scores Table
- Substitute Missing Values



- 90+ prepackaged machine learning/predictive algorithms
- Supports association, clustering, classification, regression, time series, ...
- Supports different types of data – structured, streaming and series data
- Real-time scoring for several algorithms
- Integrated with open source machine learning libraries – TensorFlow and R

# The student activity hub (SAH) will handle various needs

- Three classes of analytics

  - <u>Institutional analytics</u>: graduation rates, retention rates, enrollments, demographic, socio-economic, etc.

  - <u>Academic analytics</u>: entrance test scores, satisfactory progress, term and course grades, advising and support

  - <u>Learning analytics</u>: course engagement, assessments, clickstream, video views, discussion posts

- Other features of this architecture

  - <u>HPC type performance</u>. The system uses high-speed, in-memory techniques to handle very large data sets

  - <u>Real-time data</u>. A next generation learning environment will use real-time data for 'Fitbit' type analytics

  - <u>Personalized messaging and interactions</u>. A next generation learning environment will increasingly rely on advanced analytics to tailor degrees, content and interactions for the learner and in real-time

  - <u>Complete data integration</u>. All relevant systems are integrated in real-time or near real-time basis. Data interoperability standards will be important!

  - <u>Next-generation data science</u>. AI and ML technology can be applied to improve predictions of student success, help automate parts of summative or formative assessments, match appropriate content or other 'nudges' to learners

  - <u>Census-driven, institutional reporting and ad-hoc analysis</u>. The platform can provide both at the same time

  - <u>Highly secure</u>. Data security at multiple layers, including the database layer, the data source layer in Cognos/Tableau, within workbook creation and design, within workbook deployment in Cognos/Tableau web servers. HANA also has data anonymization built in with a real-time implementation of k-anonymity (see [http://www.sap.com/data-anonymization](http://www.sap.com/data-anonymization))

# SAH: Sample field list

| IGC Name | IGC Long Description | IGC New Status | Related Terms | IGC Permission Level | Approval person | BI Data Source | Pilot Data Source Show Field | Parent Category | Short Description | Reviewed 9/6 | Reviewed 9/13 | Reviewed 9/20 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Academic Status Code Current | Academic Status Code for end of most recently completed term or beginning of currently enrolled term. | Review2 | Academic Status Code | | Chris | Demographics | 1 | Student > Academic History >> Academic Status | | 0 | 1 | 0 |
| ACT English Score | Score received from ACT English test covering usage/mechanics and rhetorical skills. College Readiness Benchmark is 18 (2014). | Candidate | ACT | P3 | Michelle Ransom | Demographics | 1 | Student >> Admissions >> Test Scores | ACT Eng Score | 0 | 1 | 0 |
| ACT Math Score | Score received from ACT Math test covering pre-algebra, elementary algebra, intermediate algebra, plane geometry, coordinate geometry, elementary trigonometry, reasoning and problemsolving. College Readiness Benchmark is 22 (2014). | Candidate | ACT | P3 | Michelle Ransom | Demographics | 1 | Student >> Admissions >> Test Scores | | 0 | 1 | 0 |
| ACT Reading Score | Score received from ACT Reading test covering inference and understanding from the realm of prose fiction, social science, humanities, and natural science. College Readiness Benchmark is 22 (2014). | Candidate | ACT | P3 | Michelle Ransom | Demographics | 1 | Student >> Admissions >> Test Scores | ACT Read Score | 0 | 1 | 0 |
| ACT Science Score | Score received from ACT Science test covering interpretation, analysis, evaluation, reasoning, and problem-solving presented as data representation, research summary, and conflicting viewpoints. College Readiness Benchmark is 23 (2014). | Candidate | ACT | P3 | Michelle Ransom | Demographics | 1 | Student >> Admissions >> Test Scores | ACT Sci Score | 0 | 1 | 0 |
| Active Holds Academic Current Flag | Yes indicates Hold Type = AC | Candidate | Hold Type | P3 | Chris | Demographics | 1 | Student >> Student Attributes >> Holds | | 0 | 0 | 0 |
| Active Holds Administrative Current Flag | Yes indicates Hold Type = AD | Candidate | Hold Type | P3 | Chris | Demographics | 1 | Student >> Student Attributes >> Holds | | 0 | 0 | 0 |
| Active Holds Financial Current Flag | Yes indicates Hold Type = BU | Candidate | Hold Type | P3 | Chris | Demographics | 1 | Student >> Student Attributes >> Holds | | 0 | 0 | 0 |
| Age Current | Current date minus Student Birth Date | Approved | | | Chris | Demographics | 1 | Student >> Student Attributes | | 0 | 0 | 0 |
| American College Test (ACT) | The student's highest ACT (American College Testing) score. The required portion of the ACT is divided into four multiple choice subject tests: English, mathematics, reading, and science reasoning. Subject test scores range from 1 to 36. Also know as ACT Composite score. | Approved | | P3 | Michelle Ransom | Demographics | 1 | Student >> Admissions >> Test Scores | ACT Eng Score, ACT English Score | 1 | 1 | 0 |
| Attempted Units | Number of units attempted to be completed by a student. Term: The total number of units on a student record during a term. Cumulative: The total number of units on a student record. Note: Courses with a W (withdrawal notated on the transcript) count. Courses dropped without a W do not. Courses without a grade are not counted in ISIS. | Approved | | P3 | | Enrollment | 1 | Student >> Academic History >> Grades | | 1 | 1 | 0 |
| Attempted Units | Number of units attempted to be completed by a student. Term: The total number of units on a student record during a term. Cumulative: The total number of units on a student record. Note: Courses with a W (withdrawal notated on the transcript) count. Courses dropped without a W do not. Courses without a grade are not counted in ISIS. | Approved | | P3 | | Retention Detail | 1 | Student >> Academic History >> Grades | | | 1 | 0 |
| Census Date | The date of census for the term. | Review | | | | MajorMinor | 1 | Student >> Enrollment | | 0 | 0 | 0 |
| Class Department Code | A 2-4 character code in ISIS representing the Department offering the class. This code is tied to a department | Review2 | | | Chris | Enrollment | 1 | Student >> Enrollment >> Department | | 1 | 1 | 0 |
| Class Department Short Description | The Department offering the class; per Course Version data. | Review2 | | | Chris | Enrollment | 1 | Student >> Enrollment >> Department | | 1 | 1 | 0 |
| Class Division | In ISIS, Departments are tied to divisions. | Review | | | Chris | Enrollment | 1 | Student >> Enrollment >> Department | | 1 | 1 | 0 |
| Class Division ID | ID value for Division | Review | Division | | Chris | Enrollment | 1 | Student >> Enrollment >> Department | | 1 | 1 | 0 |
| Class Prior Terms Enrolled Count | Count of prior enrollments in the same course and section for the | Review | | | | Enrollment | | Student >> Enrollment | | 0 | 0 | 0 |

# "Curated views" of the data, de-identified

**Demographics**
Residency, SAT/ACT and other entrance test scores, academic status, etc.

**Enrollment**
Enrollment counts by class, departments, divisions/schools, colleges, including course grades

**Major/Minors (wide and narrow)**
Degrees, Programs, switching of majors, etc.

**Student Statistics Per Term**
Dozens of common student statistics, term-by-term for examining progression

**Class and Section Stats Per Term**
Dozens of class and section statistics, term by term for course and section planning, instructor load, etc.

**Retention (wide and narrow)**
Cohort, retention and graduation rates, etc.

**Admissions**
Applicants, Applications, Test Scores, Scholarships

**Continuing education students (Extension, other)**
Demographics, enrollment, credentials

**Learning analytics**
Learning events, grading events

**General student activities**
Activity details, Activity stats per term

| Feature_domain | Feature_Category | Feature_subcategory | Feature_ID | Feature_Name | Notes |
|---|---|---|---|---|---|
| Learning systems interactions | Session | Session | 1 | User log in | |
| Learning systems interactions | Session | Session | 2 | User log off | |
| Learning systems interactions | Session | Session | 3 | User timed out | |
| Learning systems interactions | Forums | Forum | 4 | Forum created | Created but not made available |
| Learning systems interactions | Forums | Forum | 5 | Forum posted | Made available |
| Learning systems interactions | Forums | Forum | 6 | Forum unposted | Made unavailable |
| Learning systems interactions | Forums | Forum | 7 | Forum edited | |
| Learning systems interactions | Forums | Forum | 8 | Forum deleted | |
| Learning systems interactions | Forums | Forum | 9 | Forum subscribed | |
| Learning systems interactions | Forums | Forum | 10 | Forum unsubscribed | |
| Learning systems interactions | Forums | Forum item | 11 | Forum item created | |
| Learning systems interactions | Forums | Forum item | 12 | Forum item posted | |
| Learning systems interactions | Forums | Forum item | 13 | Forum item unposted | Made unavailable |
| Learning systems interactions | Forums | Forum item | 14 | Forum item edited | |
| Learning systems interactions | Forums | Forum item | 15 | Forum item deleted | |
| Learning systems interactions | Forums | Forum item | 16 | Forum item viewed | |
| Learning systems interactions | Forums | Forum item | 17 | Forum item marked | Like, Angry, Read, Unread etc |
| Learning systems interactions | Document | Document | 18 | Document created | Created or uploaded |
| Learning systems interactions | Document | Document | 19 | Document posted | Made available |
| Learning systems interactions | Document | Document | 20 | Document edited | Re-uploaded or revised in place |
| Learning systems interactions | Document | Document | 21 | Document deleted | |
| Learning systems interactions | Document | Document | 22 | Document viewed | Document viewed or opened |
| Learning systems interactions | Assignments | Assignments | 23 | Assignment created | By instructor, created but not yet made available to students |
| Learning systems interactions | Assignments | Assignments | 24 | Assignment posted | By instructor, made available to students for access |
| Learning systems interactions | Assignments | Assignments | 25 | Assignment unposted | Made unavailable |
| Learning systems interactions | Assignments | Assignments | 26 | Assignment deactivated | By instructor, removed from access |
| Learning systems interactions | Assignments | Assignments | 27 | Assignment edited | By instructor |
| Learning systems interactions | Assignments | Assignments | 28 | Assignment deleted | By instructor |
| Learning systems interactions | Assignments | Assignments | 29 | Assignment viewed | By student |
| Learning systems interactions | Assignments | Assignments | 30 | Assignment reviewed | By instructor |
| Learning systems interactions | Assignments | Assignments | 31 | Assignment started | By student |
| Learning systems interactions | Assignments | Assignments | 32 | Assignment submitted | By student |
| Learning systems interactions | Assignments | Assignments | 33 | Assignment completed | By student |
| Learning systems interactions | Assignments | Assignments | 34 | Assignment grade created | By instructor, created, but not yet visible |
| Learning systems interactions | Assignments | Assignments | 35 | Assignment grade posted | By instructor, posted means final. There can be multiple! |
| Learning systems interactions | Assignments | Assignments | 36 | Assignment grade unposted | Made unavailable |
| Learning systems interactions | Assignments | Assignments | 37 | Assignment grade edited | By instructor, revised grade |
| Learning systems interactions | Assignments | Assignments | 38 | Assignment grade deleted | By instructor |
| Learning systems interactions | Assignments | Assignments | 39 | Assignment grade viewed | By student |
| Learning systems interactions | Assignments | Assignments | 39 | Assignment feedback created | By student or instructor |
| Learning systems interactions | Assignments | Assignments | 40 | Assignment feedback viewed | By student within the tool, not in a downloaded documenr |
| Learning systems interactions | Assignments | Assignments | 41 | Assignment feedback downloaded | e.e.g, student downloadss and assignment feedback doc |
| Learning systems interactions | Groups | Groups | 42 | Group assignment created | e.g., Instructor assigning students to a group |
| Learning systems interactions | Groups | Groups | 43 | Group assignment posted | Made available to students |
| Learning systems interactions | Groups | Groups | 44 | Group assignment unposted | Made unavailable |
| Learning systems interactions | Groups | Groups | 45 | Group assignment viewed | By the student |

# Master map of learning events

- Four level hierarchy
- At the level of granularity or lower than Caliper, xAPI
- Can map to Caliper, xAPI or future standards
- Can extend and define our learning events as needed without waiting for standards
- Can map post-hoc to standards as they evolve
- Extendible domains
  - Learning systems interactions
  - Advising interactions
  - Co-curricular interactions
  - Academic interactions
  - Advising interactions
- We are also maintaining a "Tool Hierarchy" to categorize EdTech ecosystem tools and provide a simple containership model

# Full list of curated views as of June, 2018

| | Curated view |
|---|---|
| **Student** | Student demographics |
| | Student stats per term |
| | Enrollment |
| | Retention wide |
| | Retention narrow |
| | Majors minors wide |
| | Majors minors narrow |
| | Class stats per term |
| | Section stats per term |
| | Instructor stats per term |
| | Applicants |
| | Applications |
| | Tests |
| | Scholarships |
| | *Learning events (Caliper, xAPI, LRS)* |
| | *Grading events* |
| | *Activity stats per term* |
| | *Activity details* |
| | *Student demographics CE* |
| | *Enrollment CE* |
| | *Credentials CE* |

**"In-flight" curated views**

<u>Learning analytics</u>
Learning events (Caliper, xAPI, LRS)
Grading events

<u>Student engagement</u>
Activity stats per term
Activity details

<u>Non-matriculated, extension students</u>
Student demographics
Enrollment
Credentials

# Goals: Give analysts access to anonymized views, enable real-time mobile messaging

# UC San Diego

## CAMPUS MOBILE APP

UC San Diego Campus Mobile App is a location-based mobile app that connects you to campus information such as real-time shuttles, news, events & weather.

Download on the App Store

GET IT ON Google Play

- ☼ weather and surf reports
- 🚌 location based shuttle arrival information
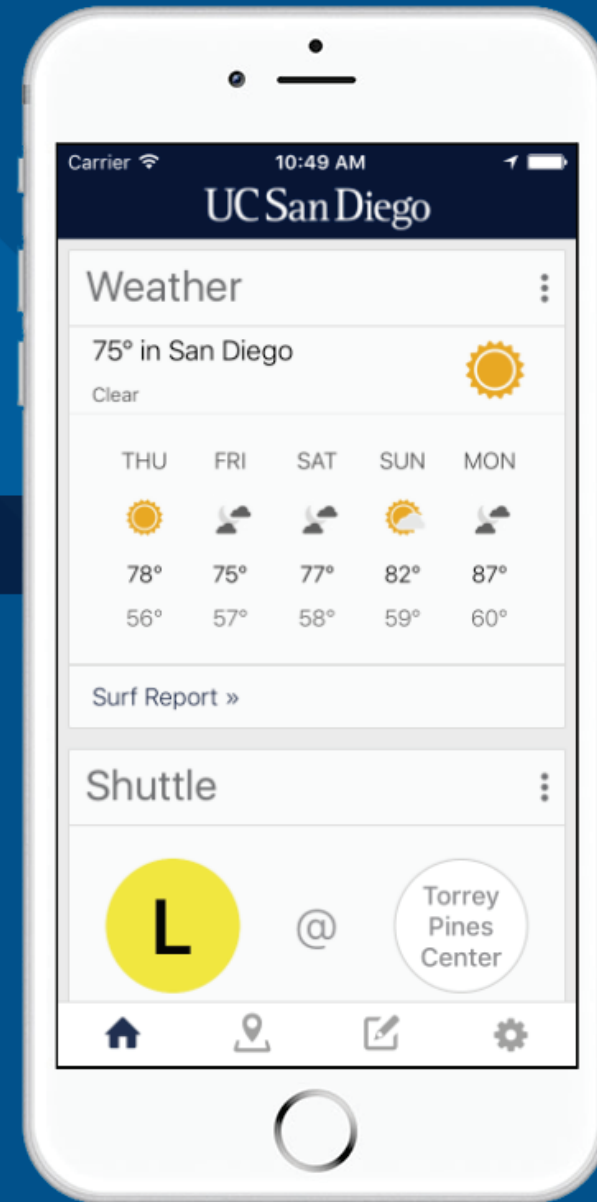- 🍴 dining menus and locations
- 📅 events updates
- ↗ links to campus services
- 📰 real-time news updates
- ➤ directions to nearby points of interest

Got feedback? Want to contribute code or design ideas for the UC San Diego Mobile app? Contact us.

### UC San Diego

Carrier 10:49 AM

**Weather** ⋮

75° in San Diego
Clear

| THU | FRI | SAT | SUN | MON |
|-----|-----|-----|-----|-----|
| 78° | 75° | 77° | 82° | 87° |
| 56° | 57° | 58° | 59° | 60° |

Surf Report »

**Shuttle** ⋮

L @ Torrey Pines Center

# SAH: Group and message builder



## Student group builder

Analyze student and learning activities to uncover trends
Filter and group students according to different attributes
Explore (and save) results in graphs and list format



## Group management

Store groups – including static and dynamic groups
Track group membership over time
Compare and analyze groups
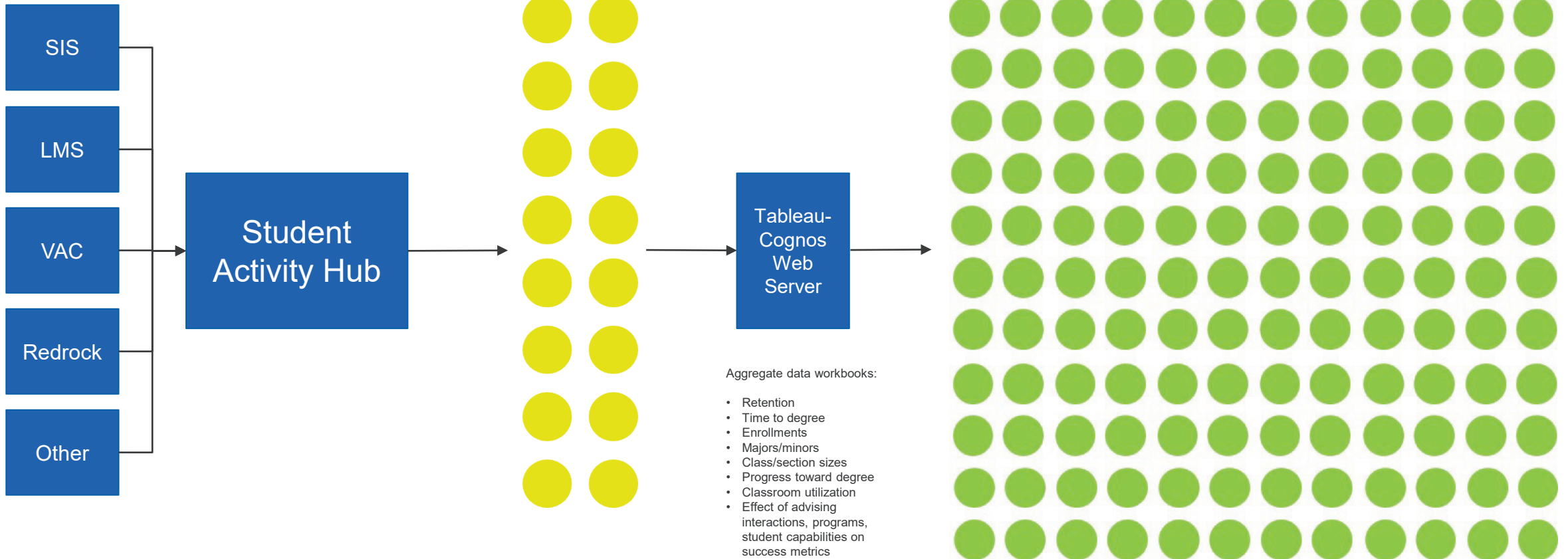Use groups as "attributes" in BI tools



## Personalized messaging

Automatically generate user-defined messages
Use message templates and embed variables
Tie message recipients to student groups

- Group builder and message builder tools interact. Group builder allows for:

- Grouping students together via any combination of fields and selection criteria (full set operations and Boolean logic)

- Changes in group membership creates events ("added to group", "removed from group") that can trigger messages, emails or workflow

- Groups also integrate with all analytics, allowing analysis to quickly compare and contrast different subpopulations of students. Subpopulations can be overlapping

- Groups are reusable and sharable and can be easily referenced within all workbooks and reports

# EXAMPLE: Student Activity Hub (SAH) Data Publishing Overview

- Legitimate educational interest only; skilled analyst
- Using Tableau Desktop, other authoring tool, secure access
- Creates dashboards, interactive analytic screens, reports
- Access to granular, de-identified data only, control small cell size if needed
- Approximately 30-40 split between central and distributed groups
- Approximately 5-8 or so publishers within primarily student service delivery offices will need identifiable data access
- Currently 70+ people have access to raw identifiable data in current DW

- Legitimate interest only; staff, faculty with secure UCSD credentials
- Accesses published workbooks via the web
- No direct data access, no identifiable data, no downloading of data
- Can manipulate the data in the workbook only to the degree the publisher allows
- Access to identifiable data, lists of students, etc. is only through the VAC or an authorized report

## Workbook viewers

## Publishers



| SIS |
| LMS |
| VAC |
| Redrock |
| Other |

**Student Activity Hub**

**Tableau-Cognos Web Server**

Aggregate data workbooks:

- Retention
- Time to degree
- Enrollments
- Majors/minors
- Class/section sizes
- Progress toward degree
- Classroom utilization
- Effect of advising interactions, programs, student capabilities on success metrics

# Hierarchy management

- Several key hierarchies need to be carefully curated so that all downstream analysis can safely aggregate and analyze data

- Hierarchy management has three components

  - **Hierarchy governance.** These are activities involving key staff who help design hierarchies, agree on publishing revised versions of hierarchies and help oversee hierarchy quality and utility. This will be included within the data and analytics governance committee

  - **Hierarchy data management.** This is a very small group of staff within the ITS Enterprise Systems team who will ensure hierarchy changes and additions are safely implemented and replicated across the subsidiary systems and activity hubs. This team can also analyze systems to determine impacts on changes to hierarchies

  - **Hierarchy data management software.** This is a software tool ITS is developing to allow the capture of key hierarchies, manage hierarchy versioning and release schedule, ensure the technical replication of hierarchy changes to subsidiary systems and allow for the mapping of different hierarchies to each other. Hierarchy mapping has interesting implications for activity hub design and use!

# Hierarchy manager tool

- This tool will enable creating and editing hierarchies, managing different versions and mapping hierarchy versions to each other independent of any enterprise system

- Enterprise systems will subscribe to one or more hierarchies as needed via the API framework where possible. Some hierarchies may be managed within an enterprise system and replicated to the hierarchy manager

- The hierarchy manager's rendering in the Activity Hubs will enable comparing and contrasting aggregate and detailed data within the curated views across different versions of a hierarchy. Example: "Show me the enrollment totals by department for the new department hierarchy for all old data." "Show me the enrollment totals by department for the old department hierarchy for all new data." "Show me those departments that have increased or decreased enrollments because of the proposed organizational change."

- Machine learning may be used to help determine new hierarchies automatically

| | Within the CoA domain | Levels | Description |
|---|---|---|---|
| | | | |
| 1 | Account | 4 | Account IDs that categorize the entry into revenue, expense, asset, liability, balance |
| 2 | Entity | 2 | Major operational unit, e.g., UC San Diego - Campus, UC San Diego - Medical Center |
| 3 | Fund | 4 | Tracks individual sources of funds |
| 4 | Department | 6 | True organizational units that have permanence and exist in org charts |
| 5 | Function | 2 | Designates the purpose of the transaction, e.g., internal, federal reporting, external reporting |
| 6 | Program | 3 | Cross-campus or system-wide program that cuts across all other hierarchies |
| 7 | Project | 3,6 | Capital and sponsored projects tracking. 3 levels in CoA, 3 more presumed needed for 6 levels total |
| 8 | Activity | 3 | General activities on campus such as Commencement, student recruiting, etc. |
| 9 | Location | 4 | Building location, vessel, etc. |
| 10 | Geolocation | 4 | Mapping coordinates for a location |
| | | | |
| | **Other relevant hierarchies** | | |
| | | | |
| 11 | Academic operations | 6 | How instructors, researchers, TAs, appointments, majors, minors, courses, degrees, programs roll up in terms of operations |
| 12 | Academic discipline | 6 | How instructors, researchers, TAs, appointments, majors, minors, courses, degrees, programs roll up in terms of discipline content |
| 13 | Employee reports to | 8 | How employees roll up to supervisors |