# Doing Research on
# Systems for Interactive Analysis
# at an Industry Lab

Leo Zhicheng Liu

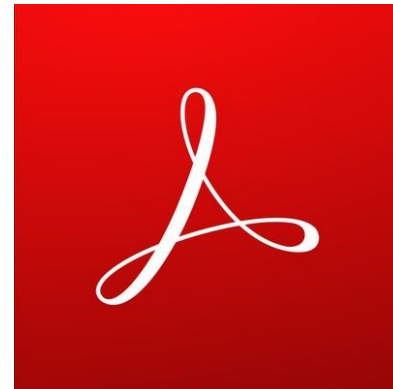Adobe Research

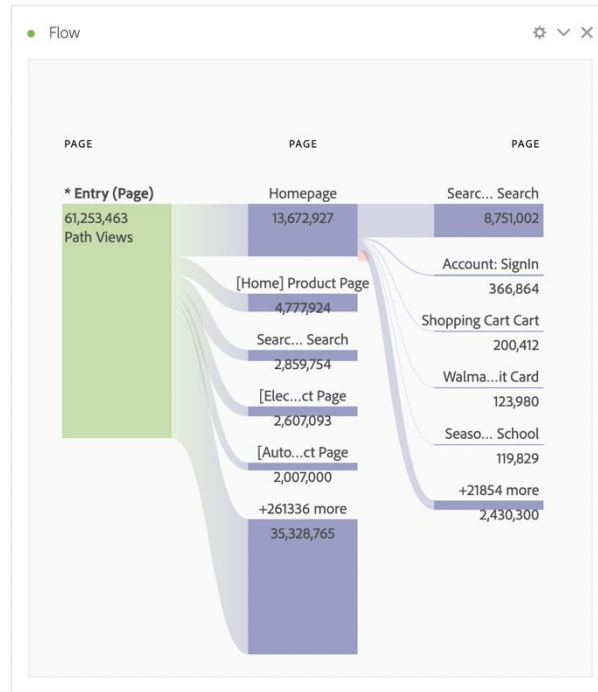**Creative**        **Marketing**        **Document**

# Adobe Analytics

"Google Analytics for Enterprises"

Tools for enterprise customers to measure and analyze customer pathing, traffic sources, content effectiveness, and video engagement.
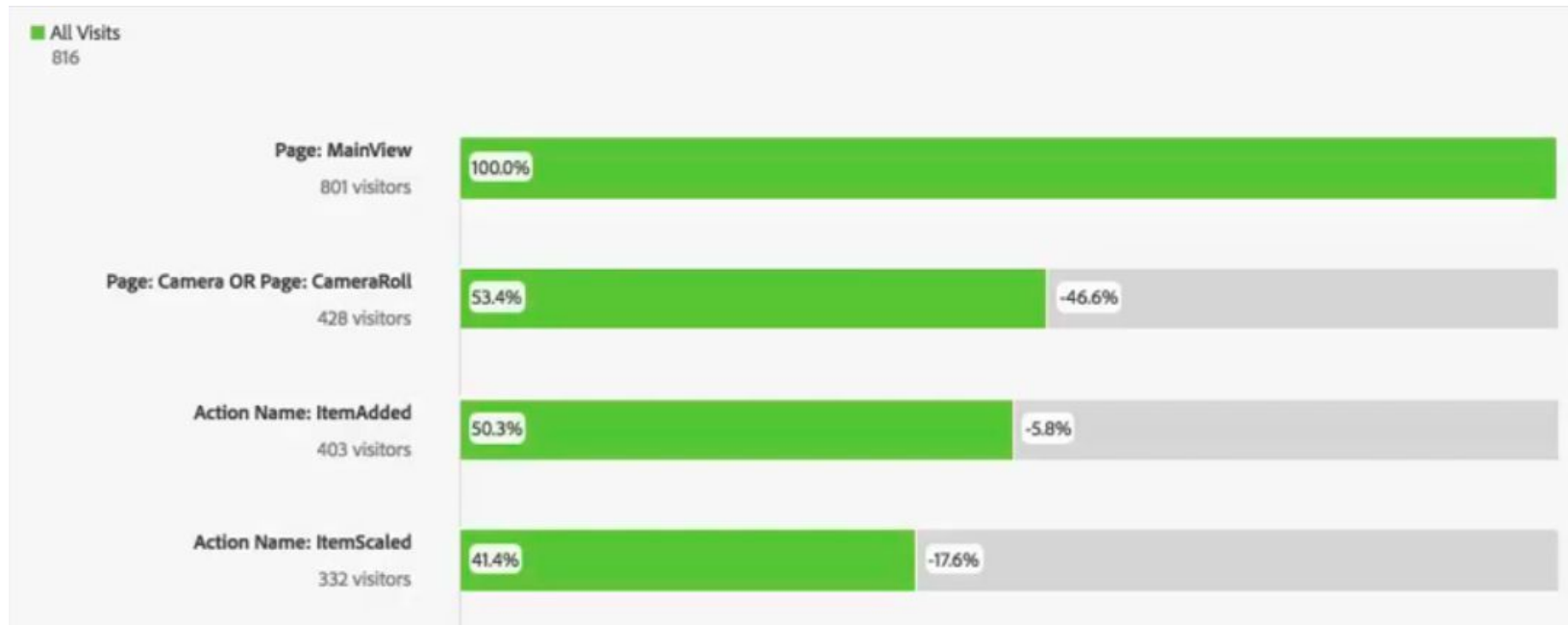
# Understanding Visitor Behavior/Journeys: Reports

Where do visitors come before/go after a specific page/site section?

# Understanding Visitor Behavior/Journeys: Reports
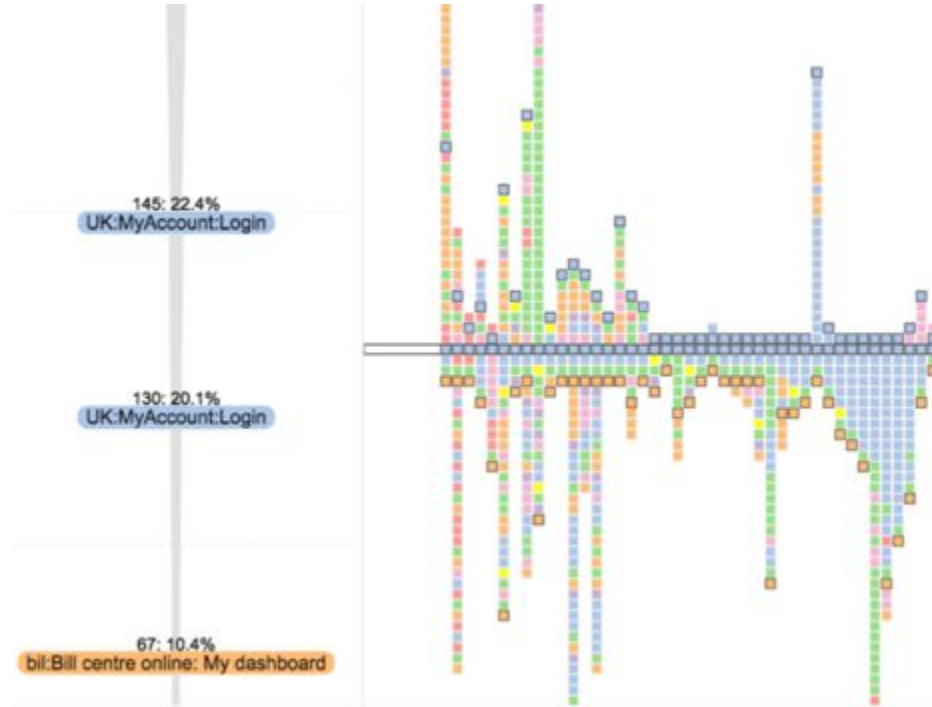
Where do visitors drop off / fall out?

# Data Management

Fast, distributed, column-oriented data store and query engine

In-house developed algebra

XML-based query language

# Interactive Analysis: Pattern Mining + Visualization



145: 22.4%
UK:MyAccount:Login

130: 20.1%
UK:MyAccount:Login

67: 10.4%
bil:Bill centre online: My dashboard

# New Mining Algorithms

# Interweaving Pattern Mining with Ad Hoc Queries

# Collaboration with Product Team

Drastic changes to the architecture/design of the datastore unlikely

Adding new queries/mining algorithms need to justify the value of approach, which is hard without real data and user feedback

Lab features to bridge the gap: test on a selected group of customers with their own data

# Collaboration with Product Team

hard, they have their priority, existing architecture cannot be easily changed to support new queries/mining algorithms that may or may not prove to be what the customers need/value

need to demonstrate the value of approach, which is hard without real data. need to show value based on small/sample data

Lab features to bridge the gap: test on a selected group of customers with their own data

Product team's support and buy-in crucial for incorporating novel ideas

# Research Collaboration

Build own datastore with smaller datasets

Work with researchers in other areas, sometimes restricted by company's talent pool

Learn from existing literature

Collaborate with academia through internships

Hire interns outside your area

The Unified Logging Infrastructure for Data Analytics at Twitter

http://vldb.org/pvldb/vol5/p1771_georgelee_vldb2012.pdf

compute histogram of event counts

construct dictionary: map each event to a symbol: more frequent events are assigned smaller code point (fewer bytes) - variable-length coding

reconstruct sessions: group by on user id and session id, 30 minute interval to delimit sessions, then encode using dictionary. A session is essentially a unicode string

only records session duration, cannot support queries about temporal gaps, lower latency, higher throughput

- summary statistics on sessions
- user sessions that match an arbitrary regular expression (event-based querying)
- funnel analytics, a funnel is a regular expression, grow the funnel substring to query number of sessions at each stage
- n-grams and pattern mining