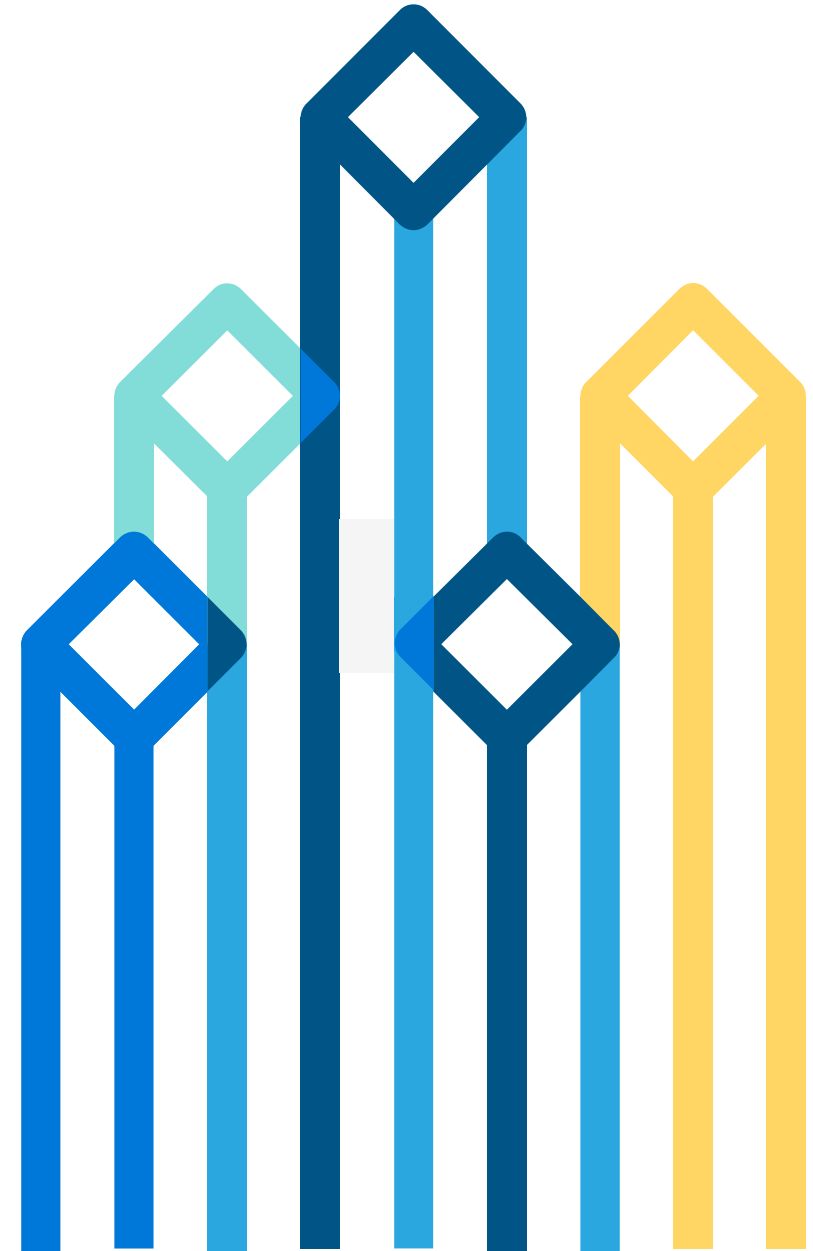# cloudera®

# Leveraging Predictive Tools to Decrease Resolution Time

Angus Klein – Vice President, Global Support

Adam Warrington – Director, Engineering

# The Value of Hadoop...

## One place for unlimited data

- All types
- More sources
- Faster, larger ingestion

## Unified, multi-framework data access

- More users
- More tools
- Faster changes

**cloudera**

# The Cloudera Value Chain



DEVELOP & PACKAGE OPEN SOURCE PROJECTS

FORM A STABLE, RELIABLE PLATFORM

THAT SUPPORTS POPULAR APPLICATIONS

TO SOLVE CUSTOMER PROBLEMS

# Problem Statement

Supporting our product is **complex**

Issues **can be related** or root cause **might not be the same**

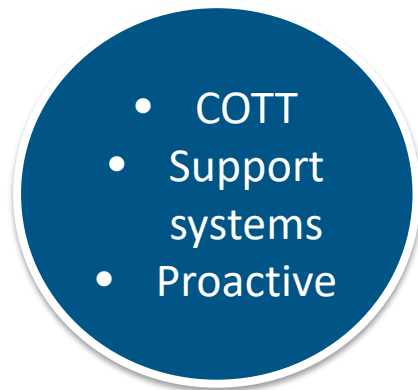Looking for ways to work globally **at scale** as company continues to grow



**cloudera**

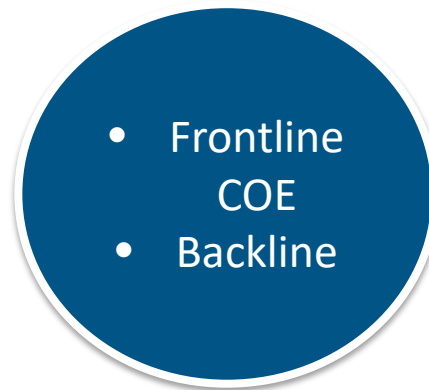# Support Team History

**Training Introduced**

2009

**Services Offered**

2010

**Support Subscription Introduced**

2011

- COTT
- Support systems
- Proactive

2012

- Frontline COE
- COE Backline

2013

**Quality Team**

2015

**cloudera**

# Enterprise grade customer support

**1st**

and only Hadoop vendor to have support certified

**1yr**

renewal cycle for certification

**12**

identified criteria areas

**100+**

service factors audited

SCP
SERVICE CAPABILITY & PERFORMANCE
- CERTIFIED -

cloudera

# Where are we?

**Global time zone coverage**

- EST (Raleigh)
- MST (Denver)
- CST (Austin)
- PST (Phoenix & Palo Alto)
- AEDT (Melbourne)
- IST (Chennai)
- JST (Tokyo)
- CST (Shanghai)
- GMT (London)
- CET (Budapest)

**cloudera**

# Structuring Predictive Support

How we integrated predictive support capabilities into our support organization

cloudera

# The Support Tools Team

Mission Statement: Build data driven software that reduces our time to solve on support cases while increasing customer satisfaction

## Key Metrics

- Decrease Time to Solve
- Increase Support Rep Throughput
- Increase Case Deflection
- Improve Support Margins

# Data Driven Support Changed the Game

Support organizations are one of the largest data drivers in any organization

Support becoming data-driven at Cloudera has been critical to establishing internal credibility at the exec table

# We dogfood

# Diagnostic Tooling – What Does This Buy Us?

**Lowers** Time to Resolution

**Improves our relationship** with Engineering

**Improves moral** of Support Engineers

**cloudera**

*Customer cases leveraging the Cloudera Diagnostic Tools demonstrated an approximate 35% drop in time to resolution.*

cloudera

# Predictive Support

Larger gains through case deflection possible with predictive support

Identifying known issues from diagnostic data

Notifying and working with the customer towards a solution to their problem

# Proactive Support

### Onboard Process

We start our partnership at the very beginning by walking you through how things work

### Predictive Validations

Powerful predictive alert system creates support tickets on behalf of our customers to help avoid known issues before they happen

### Health Check

An early warning system which looks at key indicators that represents the health of our relationship with each customer

**cloudera**

*Over 15% of support cases are deflected by Cloudera Support's predictive support system.*

# Building Predictive Support

## Step 1: Team Building

cloudera

# Team building

Create a clear and measurable mission statement

Hire 2-3 qualified engineers to prove the concept

Understand your customer – in this case, that's the supporters

# Keeping the Tools Team close

The tools team was kept close to the **full time support engineers**

Support engineers provide the **feedback loop** that allows us to build these applications

Looking for ways to work **at scale** as company continues to grow

**Tools Team**
(Debugging Tools leveraging CDH)
- CSI
- Monocle
- Clues
- Validations

**Support Systems**
- CRM System
- Community Tool
- Knowledge Base
- Support Portal
- Support Analytics

**Front Line COEs**
- Own all cases
- 24x7 Global Org
- Podded Structure
- CLR Focus

**Proactive Support**
- Customer onboarding
- Account Health Checks
- Activity Reports
- Known issue Validations

**Quality Team**
- Case Audits
- Project work to improve operations

**Back Line Support**
- Technical Escalations
- FL Mentoring
- CLR Analysis
- Supportability Focus

# Building Predictive Support

Step 2: Data Collection

cloudera

# Collect all the data



Support case interactions generates valuable support information



Troubleshooting sessions generate information about data relevant to solving a specific issue



Customer installs generate diagnostic data critical to support

# Data sources we collect

**Ingest & Consolidate**

Knowledge Base

Internal CRM Data

Support Records

Apache Mailing Lists

Community Forums

Diagnostic bundles



**OPERATIONS**
Cloudera Manager
Cloudera Director

**PROCESS, ANALYZE, SERVE**

| BATCH | STREAM | SQL | SEARCH | SDK |
|---|---|---|---|---|
| Spark, Hive, Pig MapReduce | Spark | Impala | Solr | Partners |

**UNIFIED SERVICES**

| RESOURCE MANAGEMENT | SECURITY |
|---|---|
| YARN | Sentry, RecordService |

| FILESYSTEM | RELATIONAL | NoSQL |
|---|---|---|
| HDFS | Kudu | HBase |

**STORE**

| BATCH | REAL-TIME |
|---|---|
| Sqoop | Kafka, Flume |

**INTEGRATE**

**DATA MANAGEMENT**
Cloudera Navigator
Encrypt and KeyTrustee
Optimizer

cloudera

# Data Collection Best Practices

Shadow Support Engineers to identify data and information they regularly use in a case lifecycle

Use support systems that are easy to extract data from

Store that data in a central data repository

# Game of Nodes

**50**

Nodes in cluster

**550TB**

Data size

**3TB**

New data per day

**100K**

Queries per day

**cloudera**

# Building Predictive Support

## Step 3: Visualize

cloudera

# Customer Support Interface (CSI)

**Data Ingestion**

Our internal EDH ingests 10 support specific data sources. We have access to over 500TB of data and it is growing each month.

**Data Visualization**

Our goal is to visualize all data that is useful to a support engineer in a useful way. CSI is a java web application that sits on top of the EDH

**Tools exist within CSI**

All support tools exist as a function or feature within CSI. This includes all the tools we are about to go over (e.g. Diagnostic Bundles, *Validations*, *Monocle*, and *Clues)*

# Diagnostic Bundles

**Cloudera Manager**

One of the original problems in supporting Hadoop was seeing into the customer environment. Cloudera Manager has the ability to send a snapshot using diagnostic bundles.

**What are they used for**

Support engineers are able to dive into these bundles to get a granular view of the scenario and quickly solve issues using our tools.

# Monocle

**Making specialized knowledge searchable**
Searching all of the data sources within CSI we are able to create a single index of both internal and open source knowledge for a one stop Hadoop engine.

**What is it used for**
No longer making support engineers have to "Google" for information. Our internal search platform is the most powerful Hadoop support engine for all their needs.

# Building Predictive Support

Step 4: Signature Identification

cloudera

# Closed Loop Review - Linear Process Flow

Goal: To drive supportability in the Cloudera Platform to improve the customer experience

**Step 1**

**Step 2**

**Step 3**

Support Front Line fills out "Closed Loop" data when the case is closed

Support Back Line owns analysis of Trends in customer reported issues

Input into Releases and Validations

# Issue signature creation – a collaborative process

Meet with supporters to review new spec

Develop signature code

Run against subset of last week's data

Reduce false positives and review with supporters

Release to production

# Building Predictive Support

Step 5: Delivering Predictive Support

cloudera

32

# Proactive Validations

Automatic support case creation on critical issue detection

Drive engagement with customer through known issue resolution channel

Leverage known troubleshooting mechanism and best practices

# Reactive Validations

Reactive support greatly benefits from validations

Able to show validations that might have a higher false positive rate

Able to show validations that have lower criticality, but might relate to ongoing support cases

# Basic Cluster Checklist

Run the predictive support validations at the start of a customer's contract

Getting in front of issues early saves money and increases customer satisfaction

Targeting types customers or environments that are high cost to support can improve chances of success

# cloudera
# Thank You