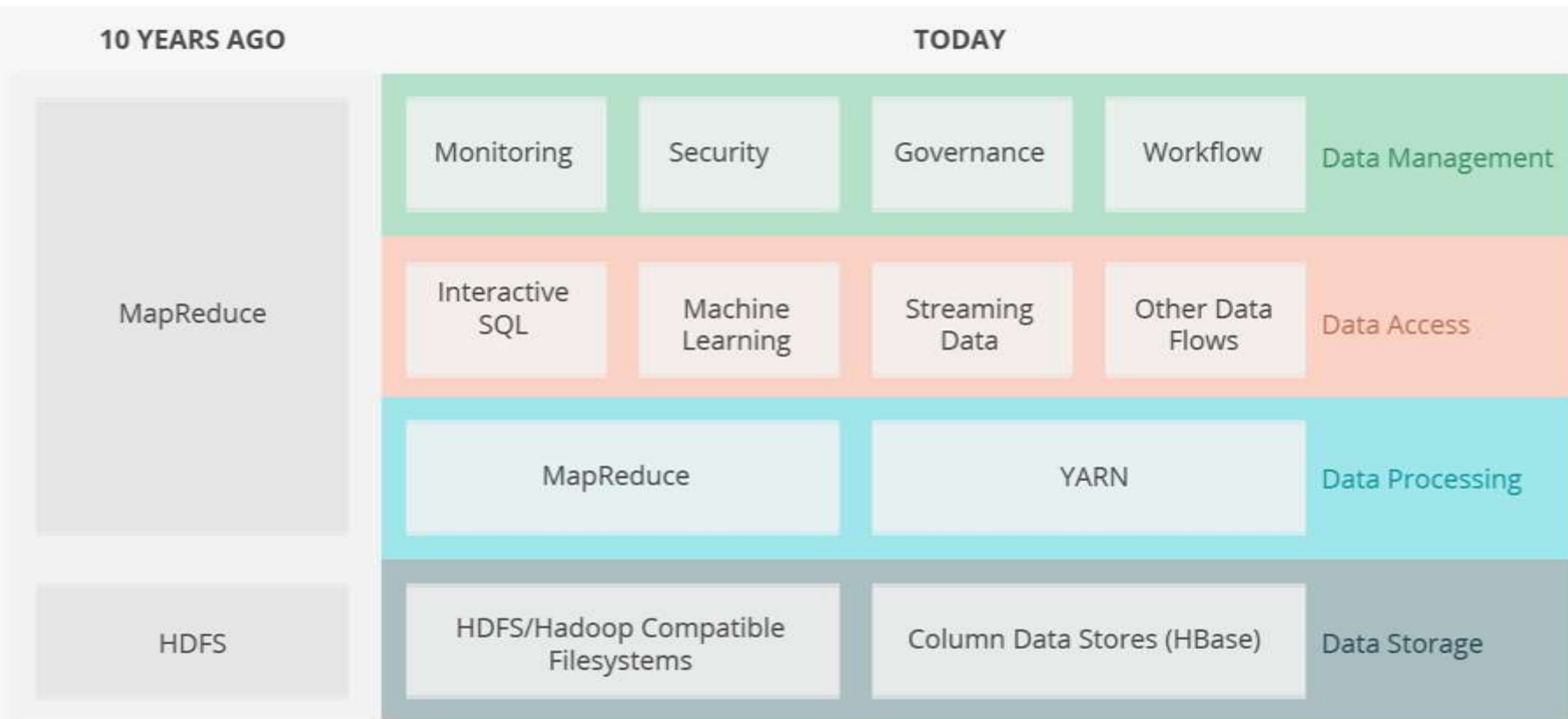


# Зачем нужен «стандартный» HADOOP?

Сергей Золотарев  
Arenadata

# Развитие экосистемы Hadoop



## Infrastructure



## Analytics



## Applications



### Cross Infrastructure /



**Open Source**



### Data Sources

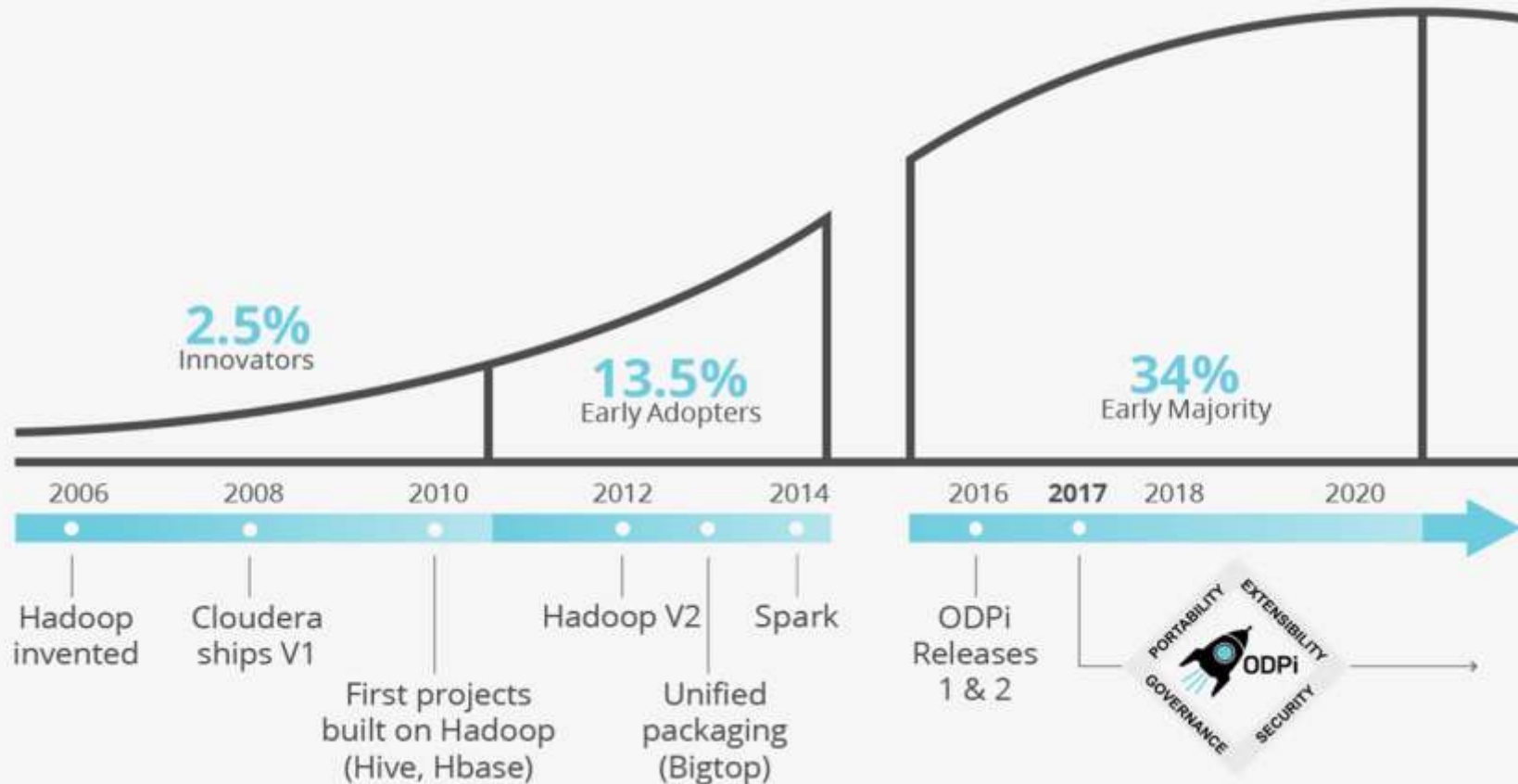


## DATA

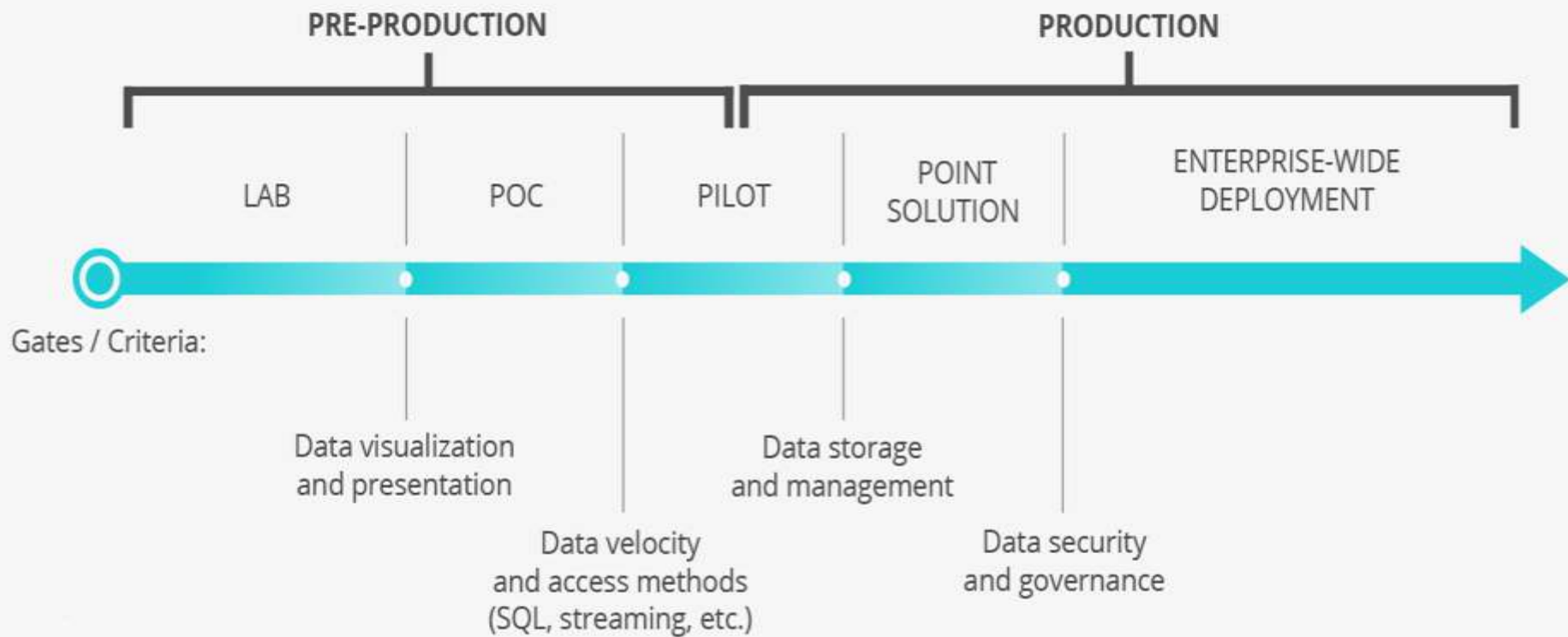




# Как корпорации «привыкают» к Hadoop



# Только 28% проектов в «Продакшн»



# Что нужно корпоративному ИТ

## SECURITY

Authenticate

Authorize

Audit

Setup Policies & Entitlements

Protect

Understand Risk Profile

## DATA GOVERNANCE

Profile

Classify

Collaborate

Understand Quality

Leverage Metadata

Provenance & Lineage

## OPERATIONS

Provision

Configure

Manage/Upgrade

Monitor

Scale

Perform



# Принципиальное отличие Пилота от Продакшн

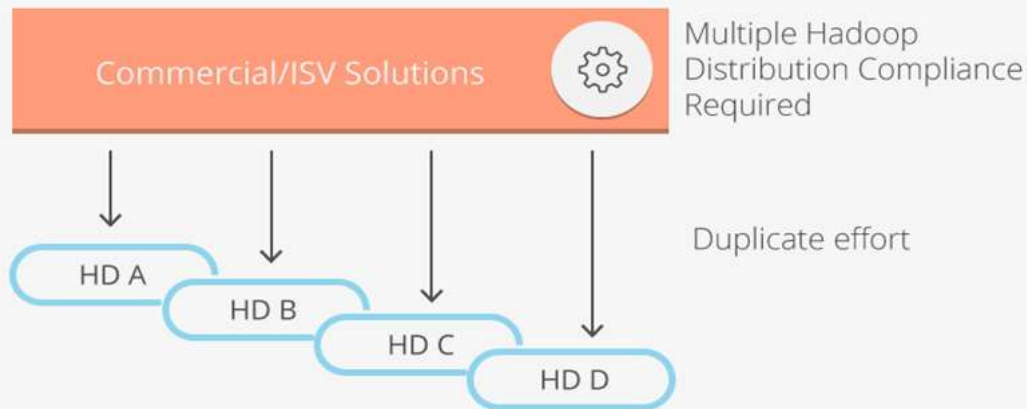
HADOOP/BIG DATA OPERATIONAL CONSIDERATIONS	LAB AND LIMITED PRODUCTION (PILOT, POC)	ENTERPRISE-WIDE PRODUCTION
<b>Deployment type</b> (single cloud, multi-cloud, on prem, single vendor, multi-vendor, etc.)	Whatever is cheapest and easiest	Tied to corporate and/or departmental procurement guidelines, and must consider SLAs, regulations, staffing, etc.
Security and data governance (data access, protection, storage, transmission, etc.)	PoCs often limit data set scope in order to keep things simple to get stakeholder buy in.	Biggest strategic gains will involve data that must adhere to strict security protections, such as customer preferences, and regulatory requirements that often require policy portability and that vary internationally.
<b>Software Lifecycle Management</b> (updating, monitoring and managing applications and platforms)	Limited usage scope means manageability and governance issues can be contained	As the number and type of uses, and the number and skill level of users, expand, the need for consistent approaches to version control, updates, application deployment, data access and other governance issues increase.
<b>Big Data / Analytics Applications</b> (ensuring uptime, extensibility and portability)	Most PoC stacks are purpose built to support a single application, so minimal need for portability and extensibility; single stack also minimizes dependency issues	With varied use-cases and tools working across shared data lakes, enterprises require a set of disciplined practices to support the objectives of each stakeholder.



# Стандартизация снижает сложность внедрения

## WITHOUT ODPI

Multi-distro certifications and regression testing increases ISV development, burden, and enterprise support costs



\*HD = Hadoop Distribution

## WITH ODPI

ODPi Interoperable Solutions



ODPi Runtime Compliant Platforms



ODPi Runtime Specification



ARENADATA

ODPi – крупнейшее мировое сообщество в области больших данных



EMC<sup>2</sup>



squid  
SOLUTIONS



NEC



IBM



**syncsort**



splunk>

**TOSHIBA**  
Leading Innovation >>>



Infosys<sup>®</sup>



**ZData** INC.



Pivotal™



vmware<sup>®</sup>

**Zettaset**  
The Leader in Big Data Security

**Xillab** [씨아랩]  
experience, idea and insight

**ARENA**DATA

**ARENA**DATA

# Вовлеченность – залог успеха



# Arenadata Hadoop : Основные факты





- ARENADATA HADOOP (ADH) – корпоративный дистрибутив платформы хранения данных Apache Hadoop
- Построен полностью на открытых компонентах проекта Apache
- Российское программное обеспечение
- Сертифицирован в 2016 году ODPI Runtime Compliant



ARENADATA



All the following Apache Hadoop platforms are [ODPi Runtime Compliant](#). This dramatically decreases engineering complexity for Big Data developers by ensuring a consistent set of base level expectations.

VENDOR	PRODUCT / VERSION	CONTACT
 <b>altiscale</b>	<a href="#">Altiscale Data Cloud 4.2</a>	Raymie Strata <rstata@altiscale.com>
	<a href="#">Arenadata Hadoop (ADH) 1.3.1</a>	Alexander Nermakov <ean@arenadata.io>
 <b>HORTONWORKS</b>	<a href="#">HDP 2.4.2</a>	Alan Gates <gates@hortonworks.com>
	<a href="#">IOP 4.2 / Biginsights 4.2</a>	Susan Malaika <malaika@us.ibm.com>

# Основные особенности

- Вся поддержка и экспертиза доступна в России и на русском языке, как удаленно так и on-site
- Разработан пакет утилит для оффлайн установки (без доступа в интернет);
- Эксплуатационная документация на русском языке
- Доступен не только в виде ПО , но и как Hadoop Appliance с полной и единой поддержкой всего программно-аппаратного комплекса от вендора
- Есть набор доступных типовых пакетных сервисов по планированию, установке и аудиту системы
- Поддержка ISV на этапе разработке решения/приложения



- ✓ HDFS
- ✓ MapReduce2
- ✓ YARN
- Tez
- ✓ Hive
- ✓ HBase
- Pig
- Sqoop
- ✓ ZooKeeper
- ✓ Flume
- ✓ Ambari Metrics
- Mahout

Actions ▾

Metrics

Heatmaps

Config History

Metric Actions ▾

Last 1 hour ▾

HDFS Disk Usage



DataNodes Live

1/1

HDFS Links

NameNode  
Secondary NameNode  
1 DataNodes

More... ▾

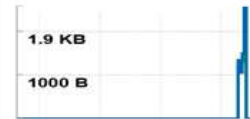
Memory Usage

18.6 GB



Network Usage

1.9 KB  
1000 B



CPU Usage



Cluster Load



NameNode Heap



NameNode RPC

0.13 ms

NameNode CPU WIO



NameNode Uptime

282.9 s

HBase Master Heap



HBase Links

HBase Master  
1 RegionServers  
Master Web UI

More... ▾

HBase Ave Load

3

HBase Master Uptime

179.2 s

ResourceManager  
HeapResourceManager  
Uptime

166.7 s

NodeManagers Live

1/1

YARN Memory



YARN Links

ResourceManager  
1 NodeManagers

More... ▾

Flume Live

1/1



**WWW.ARENADATA.IO**

**ARENADATA**