


淺談 Kubernetes 於 大數據生態系的相關開發近況



Introduction to Kubernetes Big Data
Special Interest Group (SIG)

Jazz Yao-Tsung Wang

Initiator of Taiwan Data Engineering Association
Co-Founder of Taiwan Hadoop User Group

Shared at 2017-09-21 Kubernetes 開源容器技術論壇

Hello!

I am Jazz Wang



Co-Founder of **Hadoop.TW**

Initiator of **Taiwan Data Engineering Association (TDEA)**

Hadoop Evangelist since 2008.

Open Source Promoter. System Admin (Ops).

- PAST - 11 years as a researcher in HPC field.
- 2 years (2014/03 ~ 2016/04)

Former Assistant Vice President (AVP), Product Management

You can find me at @jazzwang_tw or

<https://fb.com/groups/hadoop.tw> ,

<http://forum.hadoop.tw>

1.

前情提要

我是 K8S 初學者

只是一個 K8S Big Data SIG 觀察者

I'm a newbie to kubernetes.

I'm just an observer of K8S Big Data SIG.

關於 K8S Big Data SIG (1)

- ▷ 眾多 Kubernetes community 其中一個 SIG
 - <https://github.com/kubernetes/community>
 - Kubernetes 社群的相關 SIG 完整清單
 - <https://github.com/kubernetes/community/blob/master/sig-list.md>

Master SIG List

Name	Leads	Contact	Meetings
API Machinery	* Daniel Smith, Google * David Eads, Red Hat	* Slack * Mailing List	* Wednesdays at 18:00 UTC (biweekly)
AWS	* Justin Santa Barbara * Kris Nova, Microsoft * Chris Love * Mackenzie Burnett, Redspread	* Slack * Mailing List	* Fridays at 16:00 UTC (biweekly)
Apps	* Michelle Noorali, Microsoft * Matt Farina, Samsung SDS * Adnan Abdulhussein, Bitnami	* Slack * Mailing List	* Mondays at 16:00 UTC (weekly)
Architecture	* Brian Grant, Google * Jaice Singer DuMars, Microsoft	* Slack * Mailing List	* Thursdays at 15:30 UTC (weekly)
Auth	* Eric Chiang, CoreOS * Jordan Liggitt, Red Hat * David Eads, Red Hat	* Slack * Mailing List	* Wednesdays at 18:00 UTC (biweekly)
Autoscaling	* Marcin Wielgus, Google * Solly Ross, Red Hat	* Slack * Mailing List	* Mondays at 14:00 UTC (biweekly/triweekly)

Azure	* Jason Hansen, Microsoft * Cole Mickens, Microsoft * Jaice Singer DuMars, Microsoft	* Slack * Mailing List	* Wednesdays at 16:00 UTC (biweekly)
Big Data	* Anirudh Ramanathan, Google * Erik Erlandson, Red Hat	* Slack * Mailing List	* Wednesdays at 17:00 UTC (weekly)
CLI	* Fabiano Franz, Red Hat * Phillip Wittrock, Google * Tony Ado, Alibaba	* Slack * Mailing List	* Wednesdays at 16:00 UTC (biweekly)
Cluster Lifecycle	* Luke Marsden, Weave * Joe Beda, Heptio * Robert Bailey, Google * Lucas Kåldström, Luxas Labs (occasionally contracting for Weaveworks)	* Slack * Mailing List	* Tuesdays at 16:00 UTC (weekly)
Cluster Ops	* Rob Hirschfeld, RackN * Jaice Singer DuMars, Microsoft	* Slack * Mailing List	* Thursdays at 20:00 UTC (biweekly)
Contributor Experience	* Garrett Rodrigues, Google * Elsie Phillips, CoreOS	* Slack * Mailing List	* Wednesdays at 16:30 UTC (biweekly)

關於 K8S Big Data SIG (2)

▷ 目前的 SIG Leads

- Anirudh Ramanathan, Google
- Erik Erlandson, Red Hat (2017/8/24 剛上任)

這也是為什麼我初期只能當個觀察者～
年輕學子歡迎跳坑參加國際高手的討論

▷ 每週線上討論時間：

- 每週三 17:00 UTC = 每週四凌晨 01:00 GMT+8 (台灣時間)

▷ Slack 討論區

- <https://kubernetes.slack.com/messages/sig-big-data>
- 申請加入：<http://slack.k8s.io/>

▷ Google Group Mail List 郵件討論區：

- <https://groups.google.com/forum/#!forum/kubernetes-sig-big-data>

Big Data	<ul style="list-style-type: none">* Anirudh Ramanathan, Google* Erik Erlandson, Red Hat	<ul style="list-style-type: none">* Slack* Mailing List	<ul style="list-style-type: none">* Wednesdays at 17:00 UTC (weekly)
----------	--	--	--

關於 K8S Big Data SIG (3)

▷ 線上討論方式: Zoom 視訊/語音

- 加入方式: <https://zoom.us/my/sig.big.data>
- 每週討論 **45 分鐘 ~ 1 小時**不等

▷ 歷史會議紀錄:

- 2017/01~Now - <http://goo.gl/x5YXYS>
- 2015/10~2016/01 - <https://goo.gl/TyBB7r>

裡面有蠻多實作細節的討論
像是在哪裡卡關, 什麼元件不相容

▷ 歷史討論錄影:

- 藏在會議記錄中
- 如: 2017 九月 7 日的錄影 - <https://youtu.be/zYAyx-Wawjk>

關於 K8S Big Data SIG (4)

- ▷ 一些從 Github commit 觀察(挖)到的 SIG 歷史
 - 最早可以追溯到 **2016-05-12** 15:19 Aaron Crickenberger 提出第一次編輯
 - **2016-06-17** 10:39 sarahnovotny 說
THE BIG DATA SIG IS INDEFINITELY **SUSPENDED**,
IN FAVOR OF THE **"APPS" SIG**
 - **2017-01-30** 12:35 Anirudh Ramanathan 才又重新開啟 Big Data SIG
 - 因此歷史會議記錄是 **2017 年才比較活躍一點(目前有 74 頁)**

```
2017-01-30 12:35 Anirudh Ramanathan o Update README for sig-big-data (#307)
2016-06-17 10:39 sarahnovotny a updated to match the former wiki page
2016-05-12 15:19 Aaron Crickenberger I Add SIG pages from kubernetes.wiki as READMEs
[main] 01112e6e27f21ed081bd40210fe8ab9738904ffb - commit 17 of 17
commit 01112e6e27f21ed081bd40210fe8ab9738904ffb
Author: Aaron Crickenberger <spiffxp@gmail.com>
AuthorDate: Thu May 12 15:19:05 2016 -0700
Commit: Aaron Crickenberger <spiffxp@gmail.com>
CommitDate: Thu May 12 15:41:04 2016 -0700

    Add SIG pages from kubernetes.wiki as READMEs
---
sig-big-data/README.md | 1 +
1 file changed, 1 insertion(+)
```

關於 K8S Big Data SIG (5)

▷ Big Data SIG 的研究範疇

Covers deploying and operating big data applications (**Spark, Kafka, Hadoop, Flink, Storm**, etc) on Kubernetes. We focus on integrations with big data applications and architecting the best ways to run them on Kubernetes.

▷ Big Data SIG 的目標

- 設計相關架構, 讓大數據應用可以有效率地運行於 K8S 上

Design and architect ways to run big data applications effectively on Kubernetes

- 討論進行中的實作細節 Discuss ongoing implementation efforts

- 討論資源共享與多租戶的大數據應用

Discuss resource sharing and multi-tenancy (in the context of big data applications)

- 建議 K8S 開發有真實需求的新功能

Suggest Kubernetes features where we see a need

2.

大數據生態系整合現況

Ongoing Big Data Ecosystem Integration
with Kubernetes

ASF 目前共有 38 個大數據生態系專案



Projects Directory

[Home](#)[Committees](#)[Projects](#)[Releases](#)[Statistics](#)[Timelines](#)[About](#)

Project listings:

[By Name](#)[By Committee](#)[By Category](#)[By Programming Language](#)[By Number of Committers](#)

Projects by category:

• big-data (38):

-  [Apache Airavata](#)
-  [Apache Ambari](#)
-  [Apache Apex](#)
-  [Apache Avro](#)
-  [Apache Beam](#)
-  [Apache Bigtop](#)
-  [Apache BookKeeper](#)
-  [Apache Calcite](#)
-  [Apache CouchDB](#)
-  [Apache Crunch](#)
-  [Apache DataFu \(Incubating\)](#)
-  [Apache DirectMemory \(in the Attic\)](#)
-  [Apache Drill](#)

<https://projects.apache.org/projects.html?category#big-data>

統計方法

- 在 K8S Big Data SIG 會議記錄搜尋 38 個專案名稱當關鍵字



- Google 搜尋 38 個專案名稱 + K8S 當關鍵字

```
Tajo
Tez 12
VXQuery
Zeppelin
jazz@jazzbook ~$ for i in $(cat apache.txt | awk '{ print $3 }')
> do
> open https://www.google.com.tw/search?q=${i}+k8s
> done
```

Apache Big Data Ecosystem 整合近況一覽表

以下是 SIG 會議記錄中查得到的 Apache Big Data Project

專案	子專案	參考連結
Apache Hadoop	HDFS	<ul style="list-style-type: none"> - Data Locality Doc - https://goo.gl/zZNzwH - https://github.com/apache-spark-on-k8s/kubernetes-HDFS - https://youtu.be/DxCDxi08HWO @ Spark Summit 2017
Apache Spark	Spark Core	<ul style="list-style-type: none"> - Design Proposal - https://goo.gl/ppY28R / https://goo.gl/nyJRWi - Dynamic Allocation Proposal - https://goo.gl/QhsRaF - SPARK-18278 / Kubernetes Issue #34377 - https://github.com/apache-spark-on-k8s/spark - https://youtu.be/0xRHONrWwvU @ Spark Summit 2017
Apache Zeppelin		<ul style="list-style-type: none"> - 搭著 Spark 順風車 https://github.com/kubernetes/kubernetes/tree/master/examples/spark
Apache Storm		https://github.com/kubernetes/kubernetes/tree/master/examples/storm
Apache Cassandra		<ul style="list-style-type: none"> - https://kubernetes.io/docs/tutorials/stateful-application/cassandra/ - https://github.com/kubernetes/examples/tree/master/cassandra
Apache Kafka		<ul style="list-style-type: none"> - https://github.com/kubernetes/contrib/tree/master/statefulsets/kafka
Apache Airflow		<ul style="list-style-type: none"> - Roadmap - https://goo.gl/BpM4jq



Apache Big Data Ecosystem 整合近況一覽表

以下是 Google Apache Big Data Project + K8S 找到的

專案	子專案	參考連結
Apache Hadoop	YARN	YARN、Mesos、K8S 的定位很接近, 目前看到 YARN on K8S 的實作 - https://github.com/Comcast/kube-yarn Docker & Kubernetes on Apache Hadoop YARN https://hortonworks.com/blog/docker-kubernetes-apache-hadoop-yarn/
Apache Ambari		- https://github.com/davidstack/docker-ambari
Apache Beam		- https://github.com/apache/beam/tree/master/.test-infra/kubernetes
Apache Bookkeeper		- http://bookkeeper.apache.org/docs/latest/deployment/kubernetes/ - Bookkeeper issue #337 - https://github.com/fcunydistributedlog-on-k8s/blob/master/bookkeeper.statefulset.yaml
Apache CouchDB		CouchDB 2.0 in Kubernetes - https://gist.github.com/kocolosk/d4bed1a993c0c506b1e58274352b30df
Apache Drill		https://hub.docker.com/r/jowanza/apache-drill/ https://josep2.github.io/Jathena/

Apache Big Data Ecosystem 整合近況一覽表

以下是 Google Apache Big Data Project + K8S 找到的

專案	參考連結
Apache Flink	<ul style="list-style-type: none">- https://ci.apache.org/projects/flink/flink-docs-release-1.3/setup/kubernetes.html- 官方有 docker image https://hub.docker.com/_/flink/- FLINK-5966 / kubernetes issues #15817
Apache Flume	<p>在 kubernetes 上使用 Flume TAILDIR 收集日誌到 HDFS 上</p> <ul style="list-style-type: none">- https://ieevee.com/tech/2017/05/11/flume.html- https://github.com/vishnudxb/kube-ignite
Apache Ignite	<p>Kubernetes and Apache® Ignite™ Deployment on AWS</p> <ul style="list-style-type: none">- https://www.gridgain.com/resources/blog/kubernetes-and-apacher-ignitetm-deployment-aws
Apache Kafka	https://github.com/kubernetes/charts/tree/master/incubator/kafka

Spark + Zeppelin on Kubernetes



kubernetes

An open source system for automating deployment, scaling, and operations of applications.

[Learn about Kubernetes](#)

Wednesday, March 30, 2016

Using Spark and Zeppelin to process big data on Kubernetes 1.2

Editor's note: this is the fifth post in a [series of in-depth posts](#) on what's new in Kubernetes 1.2

With big data usage growing exponentially, many Kubernetes customers have expressed interest in running [Apache Spark](#) on their Kubernetes clusters to take advantage of the portability and flexibility of containers. Fortunately, with Kubernetes 1.2, you can now have a platform that runs Spark and Zeppelin, and your other applications side-by-side.

Why Zeppelin?

[Apache Zeppelin](#) is a web-based notebook that enables interactive data analytics. As one of its backends, Zeppelin connects to Spark. Zeppelin allows the user to interact with the Spark cluster in a simple way, without having to deal with a command-line interpreter or a Scala compiler.

Subscribe To Blog

Posts

Comments



[@Kubernetesio](#)



[View on GitHub](#)



[#kubernetes-users](#)



[Stack Overflow](#)



[Download Kubernetes](#)

Blog Archive

<http://blog.kubernetes.io/2016/03/using-Spark-and-Zeppelin-to-process-Big-Data-on-Kubernetes.html>

Spark 2.2 已將 K8S 列為實驗叢集管理

[Overview](#)[Programming Guides ▾](#)[API Docs ▾](#)[Deploying ▾](#)[More ▾](#)

Cluster Manager Types

The system currently supports three cluster managers:

- [Standalone](#) – a simple cluster manager included with Spark that makes it easy to set up a cluster.
- [Apache Mesos](#) – a general cluster manager that can also run Hadoop MapReduce and service applications.
- ~~[Hadoop YARN](#) – the resource manager in Hadoop 2.~~
- [Kubernetes \(experimental\)](#) – In addition to the above, there is experimental support for Kubernetes. Kubernetes is an open-source platform for providing container-centric infrastructure. Kubernetes support is being actively developed in an [apache-spark-on-k8s](#) Github organization. For documentation, refer to that project's README.

Submitting Applications

Applications can be submitted to a cluster of any type using the `spark-submit` script. The [application submission guide](#) describes how to do this.

Monitoring

Each driver program has a web UI, typically on port 4040, that displays information about running tasks, executors, and storage usage. Simply go to `http://<driver-node>:4040` in a web browser to access this UI. The [monitoring guide](#) also describes other monitoring options.

<http://spark.apache.org/docs/latest/cluster-overview.html>

< 插播 >

工商服務時間

台灣資料工程協會

Taiwan Data Engineering Association

Let's Play with Data Together !!

台灣資料工程協會公開徵求會員

台灣資料工程協會個人會員線上申請 表單

本會經內政部 106 年 7 月 10 日台內團字第 1061401953 號函准籌組,並 成立籌備會,茲公開徵求會員。

一、本會宗旨:

本會為依法設立、非以營利為目的之社會團體,以結合理論與實務推廣資料工程技術與應用,增進國內與國際社群交流,培育專業人才回饋社會為宗旨。

二、入會資格:

個人會員：年滿二十歲，對資料工程及其相關應用有興趣，且認同本會宗旨者，填具入會申請書，經理事會通過並繳納會費者，為個人會員。

學生會員：年滿二十歲，就讀於各大專院校在校學生，填具入會申請書，經理事會審查通過者，得加入本會為學生會員。

永久會員：個人會員自願一次繳交永久會費者，為永久會員。

榮譽會員：對本會事業有卓越貢獻，經理事會通過者，為榮譽會員。

贊助會員：對本會熱心贊助或捐助款項達一定金額，經理事會通過者，為贊助會員。

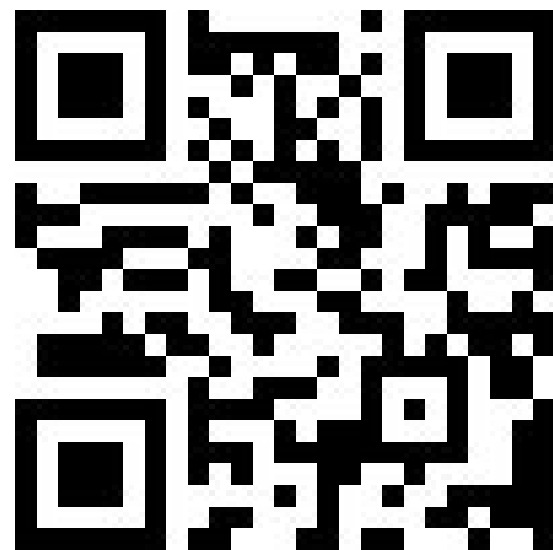
團體會員：認同本會宗旨之機構或團體，填具入會申請書，經理事會通過，並繳納會費後成為本會團體會員。每名團體會員得推派兩位代表，以行使會員權利。

三、籌備期間申請入會之截止日期:即日起至 106 年 9 月 25 日止。

四、籌備會聯絡信箱：info@dataengineering.tw。

個人會員線上申請表單

: <https://goo.gl/2z9BGK>



1st Apache Contributor Hackathon

- ▷ 台灣不該只是技術使用者(接收), 更該晉級技術開發者(供給)
- ▷ **We have 6 Apache Committer in Taiwan !!**
 - 蔡東邦 - Apache Spark Committer (台大/成大物理)
 - 陳恩平 - Apache Mesos Committer
 - 葉祐欣 - Apache BigTop Committer (現任 BigTop Project Chair, 成大資管)
 - 莊偉超 - Apache Hadoop Committer (交大)
 - 戴資力 - Apache Flink Committer (成大)
 - 蔡嘉平 - Apache HBase Committer (成大資工)



第一屆 Apache Contributor 育成賽

<https://goo.gl/6JBDzD>



3.

結語：從 SIG 學到的事情

Lessons Learned from K&S Big Data SIG

結語：我從 SIG 學到的事情

- ▷ 台灣要國際化，K8S SIG 提供跨時區協同作業的良好範例
 - Zoom 視訊 / Slack / Google Docs 會議記錄 / YouTube 錄影 / Github 版控
- ▷ 建議多參與國際自由軟體的 SIG 可以擴展自己的視野
 - 跟 Google, Redhat 等大型軟體公司的程式高手交手的機會
 - 看 Apache Software Foundation 的 JIRA 跟 Github 的 Issue 學習軟體工程 / CI/CD 的 Best Practice
- ▷ 進化論：
 - 使用者 -> 參與 SIG 的開發討論 -> 成為開發者