# Course Introduction
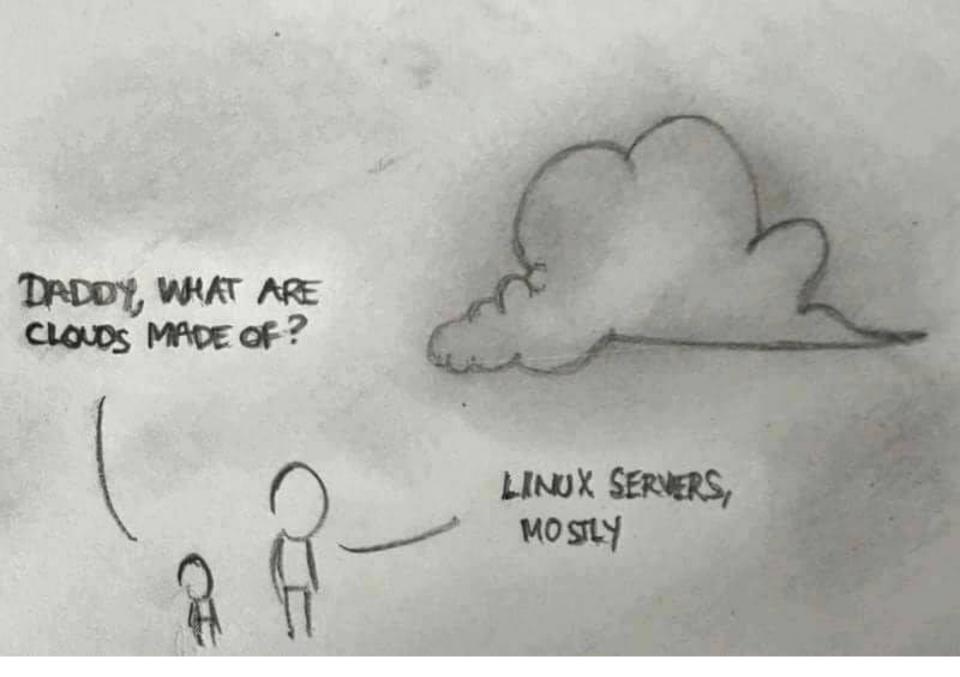
# Big Data Engineering in the Cloud

Dec 2017

# Introduction

- Aims
- Pre-requisites
- Contents
- Objectives
- Resources
- Rules of Engagement
- Introductions

# Aims

- Understanding principles of Big Data
- Theoretical background and origins
- Practical experience of modern big data processing systems technologies
- Architecture and design
- Wider context

# Pre-requisites

**Covered by the Pre-Study Guide**

- **Command line** tooling and Unix commands
- Some **Python programming** and **text editors**

# Format

- A mixture of lectures and practical labs

- Lectures aim to provide the wider context and background
  - Independent of specific technologies

- Labs are based on specific technologies
  - Designed to demonstrate the principles

# Lab model

- Local Virtual Machine
  - Ubuntu
  - Pre-installed big data software
    - E.g. Apache Spark, Cassandra, Python
- Amazon Web Services
  - Virtual machines in the cloud

# Contents

- Big Data motivation and overview
- Using Python for Data Analysis
- Map Reduce and Directed Acyclic Graphs
- Apache Spark
- Spark and SQL
- Theory of scaling
- Running Spark on Amazon
- Introduction to NoSQL databases
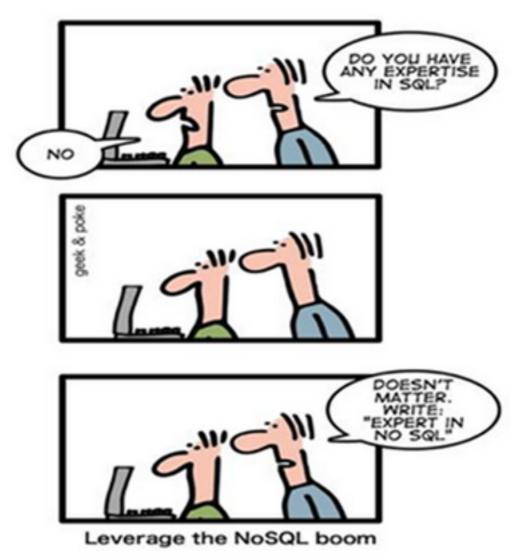- Introduction to Machine Learning

# Practicals

- Python Data Analysis
- Spark, SparkSQL
- Spark on Amazon
- Cassandra and NoSQL
- Machine Learning libraries
- Visualisation

# Improve your CV?

# Beyond the scope of this course

- Detailed Data Science techniques
- Understanding **all** of Hadoop, Spark, HDFS, Machine Learning

# Rules of Engagement

- ***Ask questions as we go along***
  - We will "park" any that are better answered later
  - Don't wait till the end to ask or raise concerns
  - If you don't ask we can't help you

# There ~~might~~ will be bugs!

- Please help out:
  - Please create new issues on the Github repository
  - https://github.com/pzfreo/big/issues/new

# Paul Fremantle

- CTO and Co-Founder of WSO2
- Previously Senior Technical Staff Member, IBM WebSphere architecture
- Visiting Lecturer, Oxford University
- VP, Apache Synapse and Member of Apache
- PhD in Computing (2017)
  - IoT security and privacy

# David Johnson

- Senior Researcher, e-Science, Oxford University

- Founding Member, Data Science Institute, Imperial

- PhD, Reading University, 2010

- Awarded University of Oxford Teaching Award, 2016

# You?

# Approximate Schedule

- Weds
  - Introductions
  - Overview and Motivation
  - Data Analysis with Python and Pandas
  - Map Reduce
  - Apache Spark

- Thursday
  - SQL
  - Theoretical background on scaling systems
  - Scaling Spark on AWS
  - Visualisation

- Friday
  - Introduction to Machine Learning
  - Realtime systems
  - Architecting big data systems
  - Completion of labs

# Let's get started