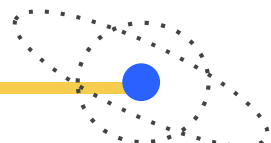
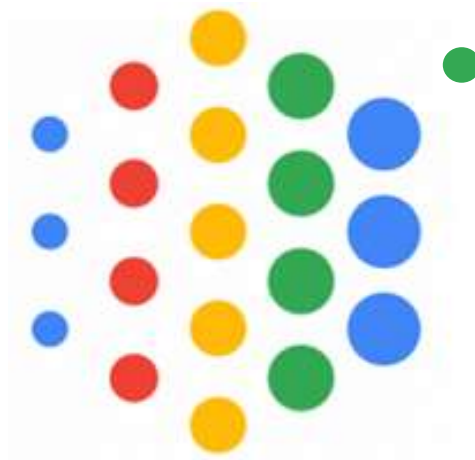


Intro to Machine Learning

<Vincent Tatan>
vintatan@google.com



Wait, who am I?

Tak Kenal maka Tak Sayang,
Tak Sayang, maka Tak Tanya
Tak Tanya, maka Tak Tahu
-- Vincent Tatan --





Meet Vincent

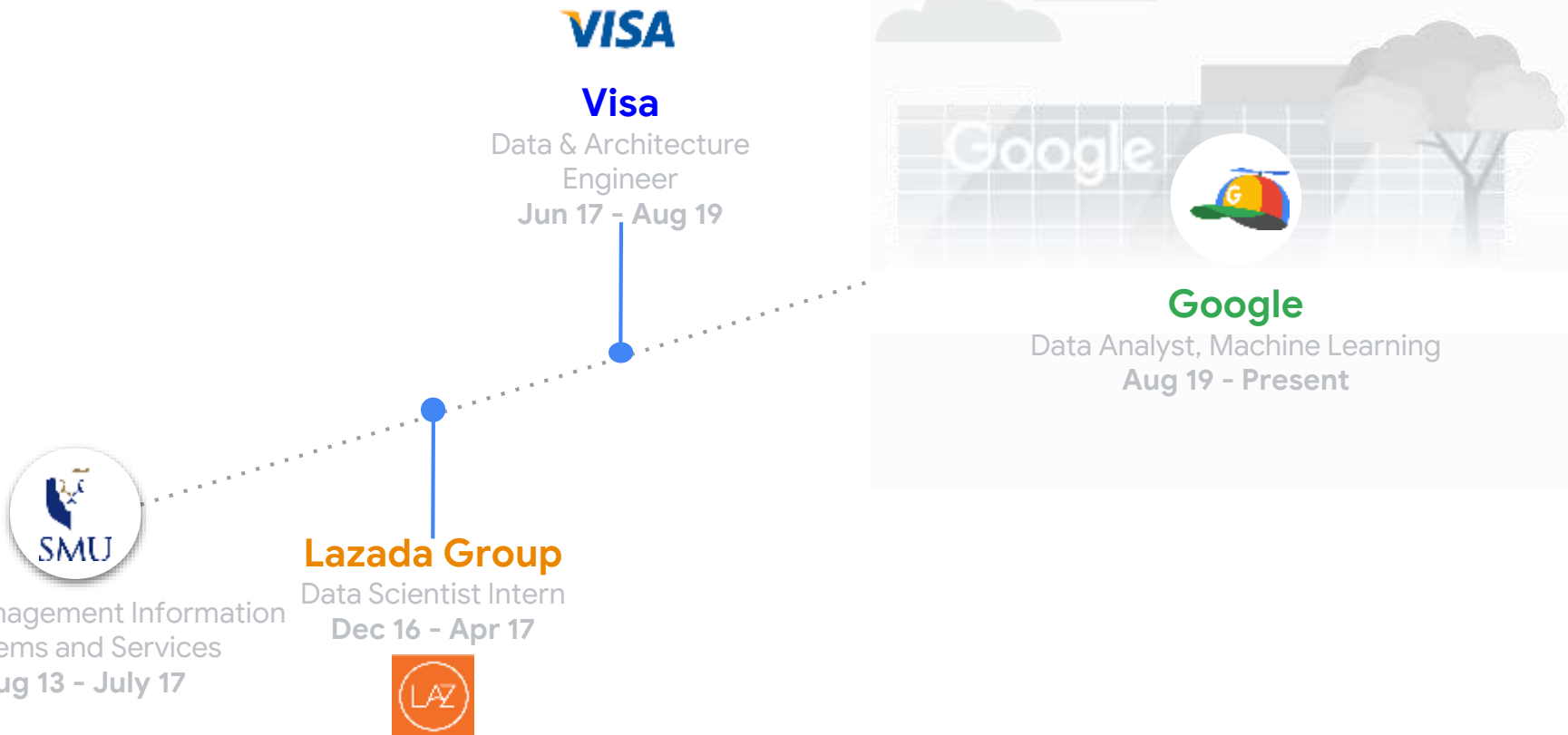
Safe Browsing Analyst (Machine Learning)
Google Trust & Safety

Medium : towardsdatascience.com/@vincentkernn

Linkedin : linkedin.com/in/vincenttatan/

Data Podcast: <https://datacast.simplecast.com/>

Path to Google



Google: Trust and Safety



To **prevent phishing @
scale** with **Data Analytics
and ML**

So, what do I do at Google?

SAFETY AND SECURITY

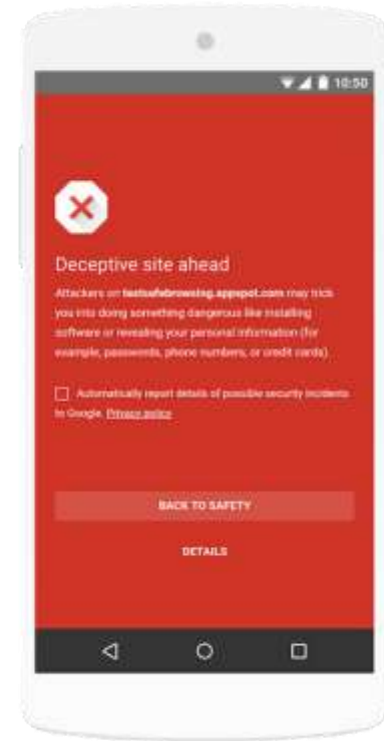
The geeky detective-work that protects you online, automatically



Trust and Safety

Protect more than 3 billion devices worldwide

1. **Google notifies your browsers** to prevent **phishing** and **malware**.
2. Using machine learning-based detection, we contributed to **99.9% accuracy in spam detection**
3. So if you see this, **beware!**



What is machine learning?



Machine Learning

Computational methods using **experience** to improve
performance

Machine Learning

Using **Computer** and **data** to **achieve objective**

- **Computer** → Algorithm, complexity analysis, theoretical guarantees.
- **Data analysis** → Statistics, probability
- **Achieve Objective** → Understanding the problem, simulation, evaluation, etc

Machine Learning



What society
thinks I do



What stock holders
think I do



What my manager
thinks I do



What product teams
think I do



What I think I do



What I really do

Supervised vs Unsupervised



Unsupervised Learning : No labeled data. Finding patterns/insights



Unlabeled
Training Data

The diagram consists of a single orange oval containing the text 'Unlabeled Training Data'.

Supervised Learning: Most common learning scenarios

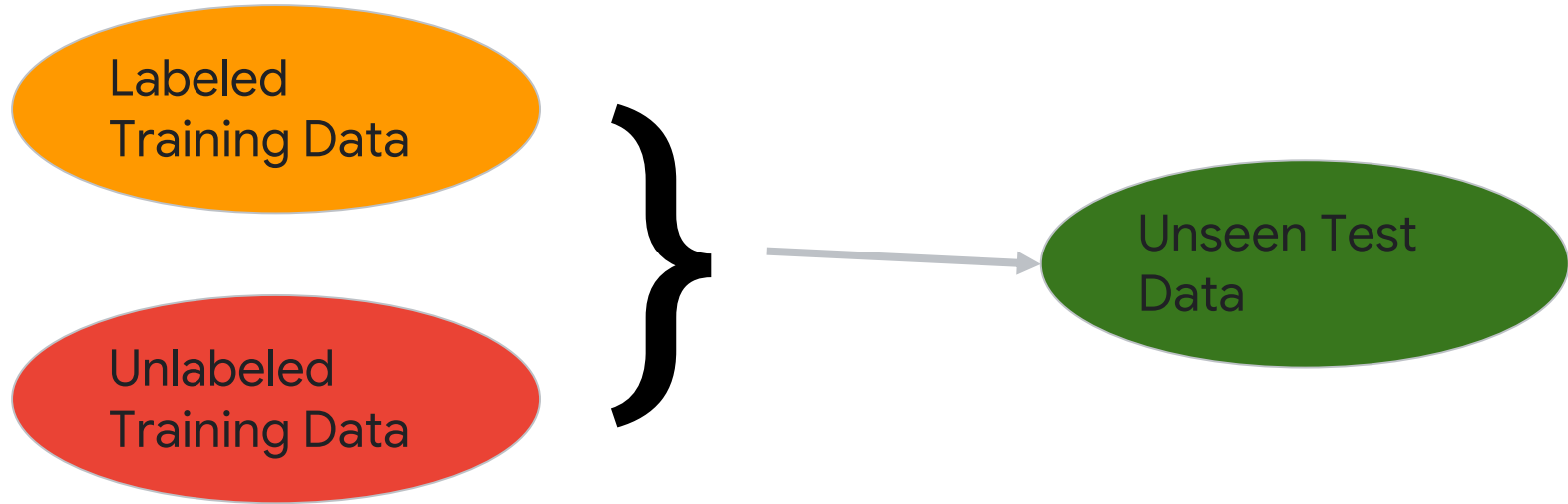


Labeled
Training Data

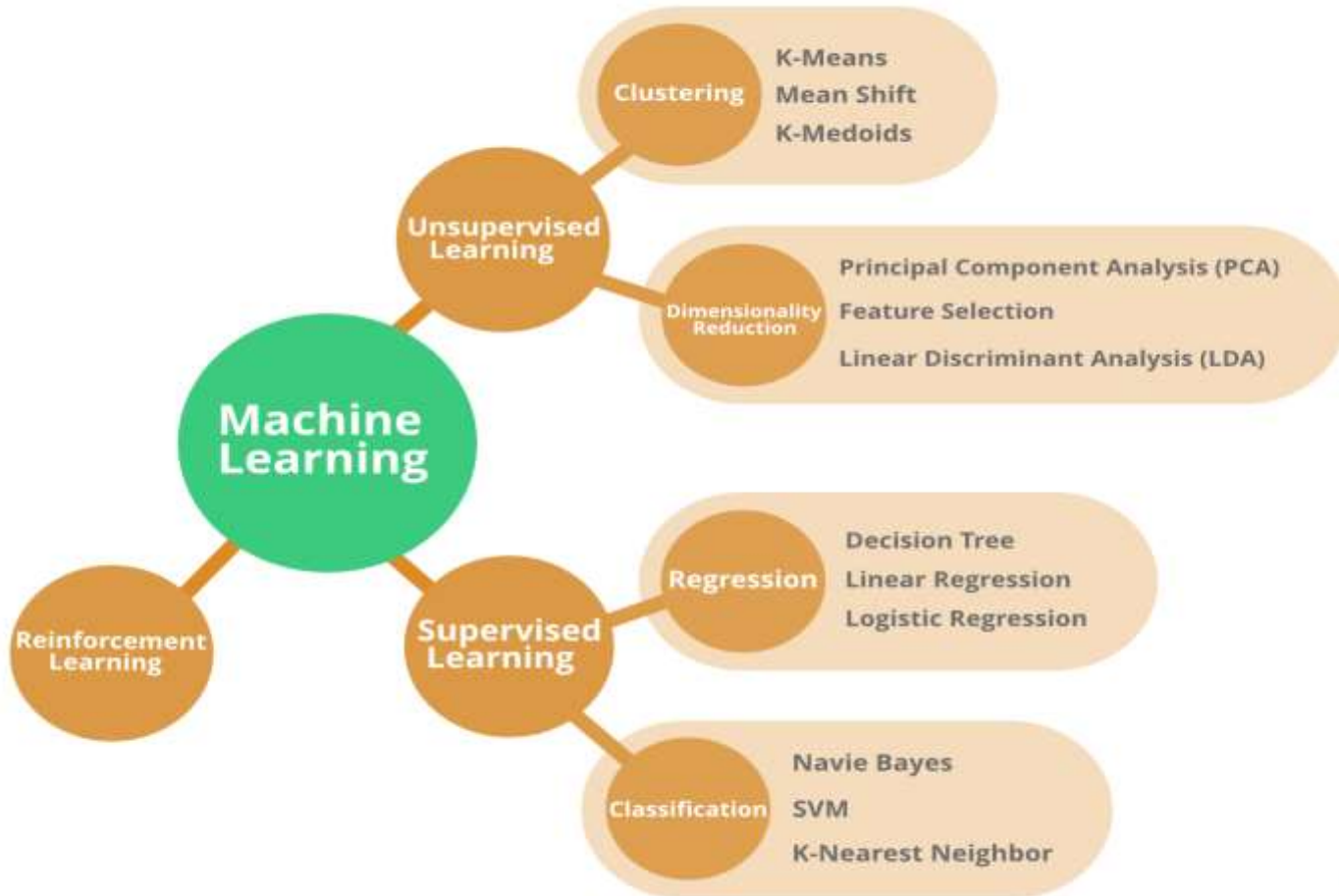
The diagram shows an orange oval on the left containing the text 'Labeled Training Data'. A horizontal grey arrow points from the right side of this oval to a green oval on the right. The green oval contains the text 'Unseen Test Data'.

Unseen
Test Data

Semi Supervised Learning : With labeled and unlabeled training data



Why? Training data might imply same distributions.



Real World Impact

What are the applications of AI and ML?



Important Resources (Teachable Machine)

<https://teachablemachine.withgoogle.com/train/image>

The screenshot displays the Teachable Machine web interface. On the left, three classes are defined: 'Rock' with 34 image samples, 'Scissors' with 26 image samples, and 'Paper' with 29 image samples. Each class has 'Webcam' and 'Upload' buttons. A 'Training' panel in the center shows 'Model Trained' and an 'Advanced' dropdown. On the right, the 'Preview' panel shows a live webcam feed of a person and the model's output probabilities: 'Rock' at 0%, 'Scissors' at 100%, and 'Paper' at 100%.

Teachable Machine

Rock 34 Image Samples

Webcam Upload

Scissors 26 Image Samples

Webcam Upload

Paper 29 Image Samples

Webcam Upload

Add a class

Training

Model Trained

Advanced

Preview Export Model

Input: ON Webcam

Switch Webcam

Output

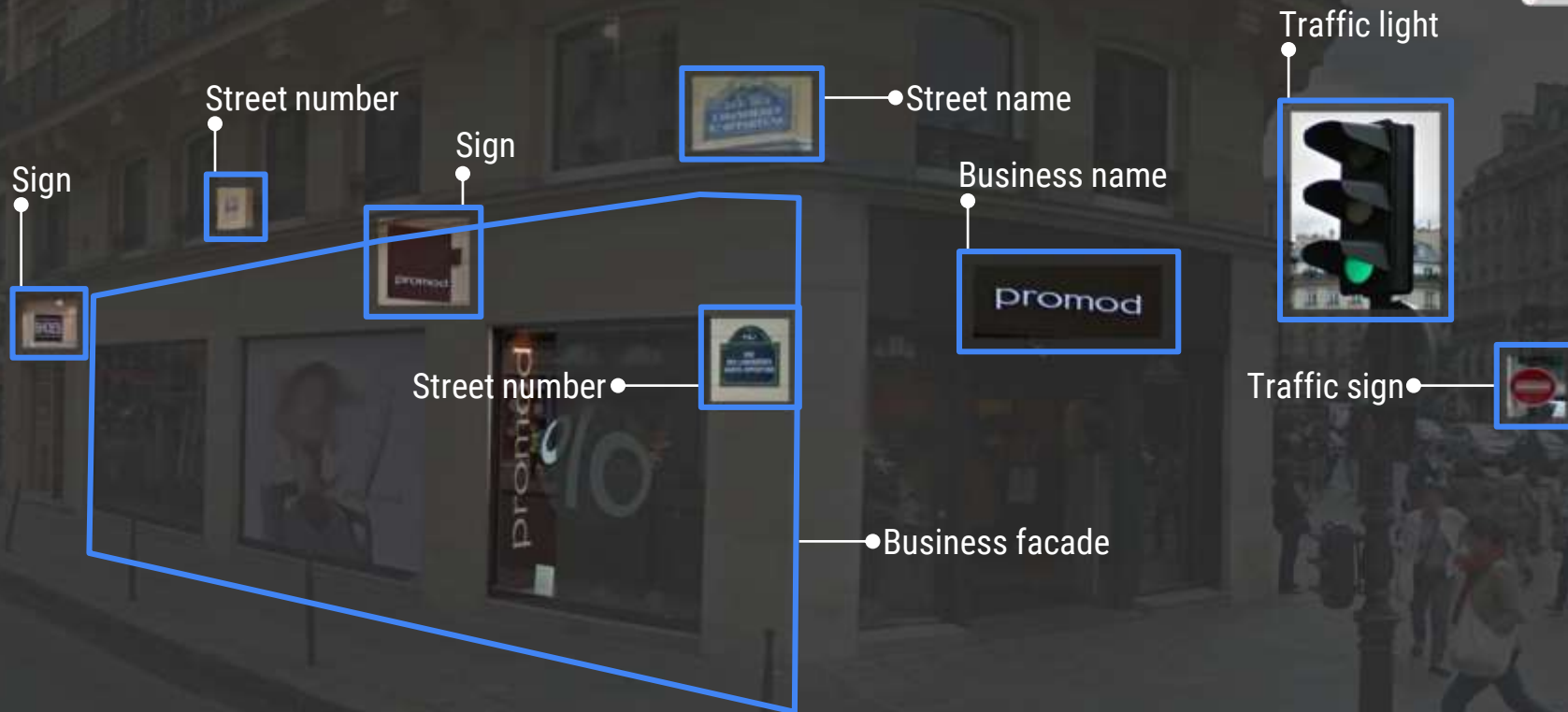
Rock 0%

Scissors 100%

Paper 100%



Street View





Google Translate

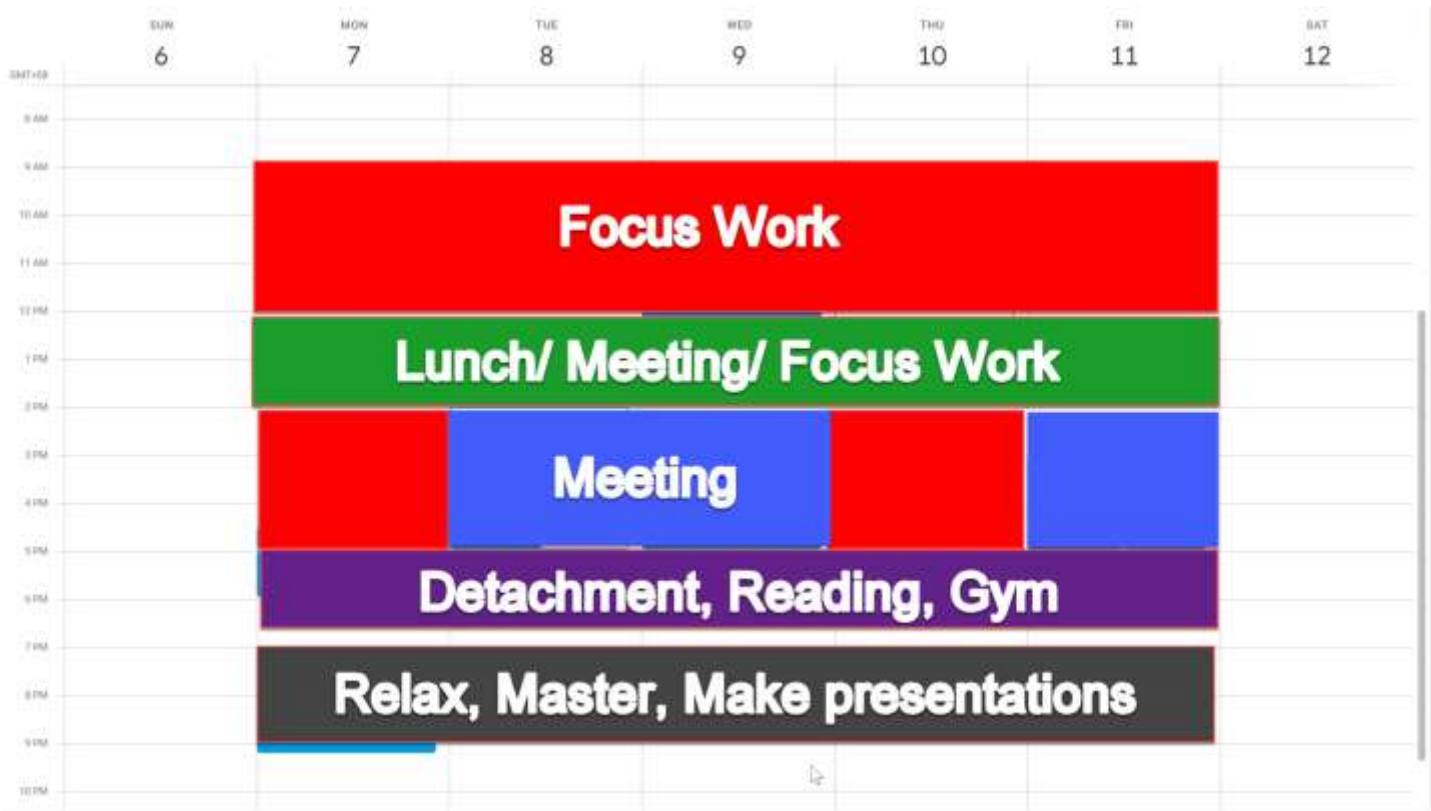


What do I do mostly?



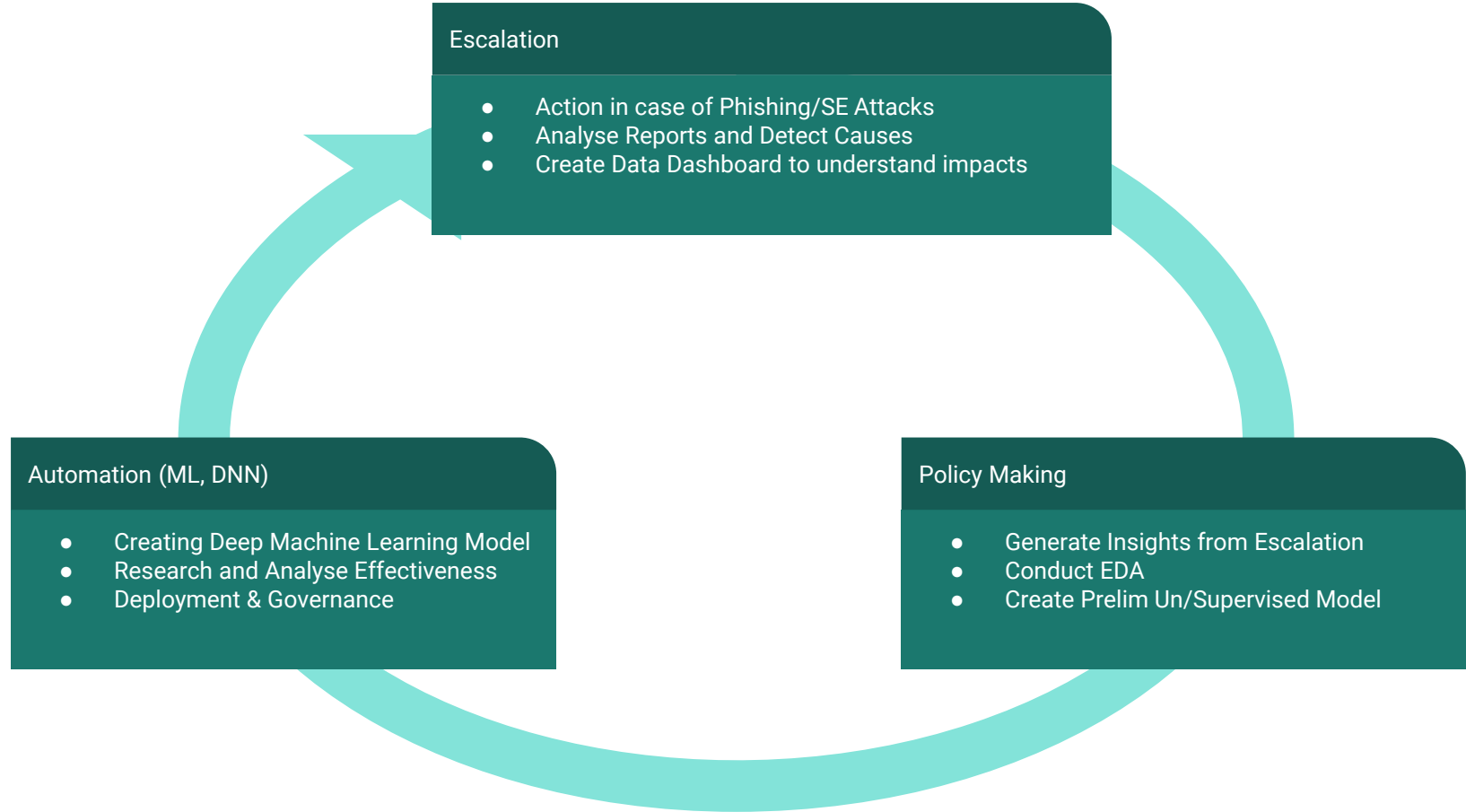
My Life at Google...

Proprietary + Confidential



Focus Work: The cycle of Data Project

Proprietary + Confidential



ML Pipeline



Frame the Problem

Proprietary + Confidential

What is your goal?

Who are your stakeholders?

How do you add value to them?



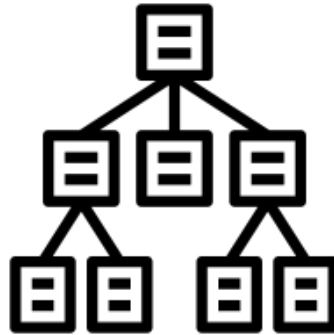
ML Pipeline

Proprietary + Confidential

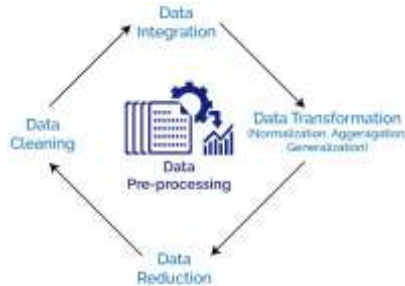
Data Collection +
Preprocessing



Model Training
and Evaluation



Machine Learning
Operations (MLOps)



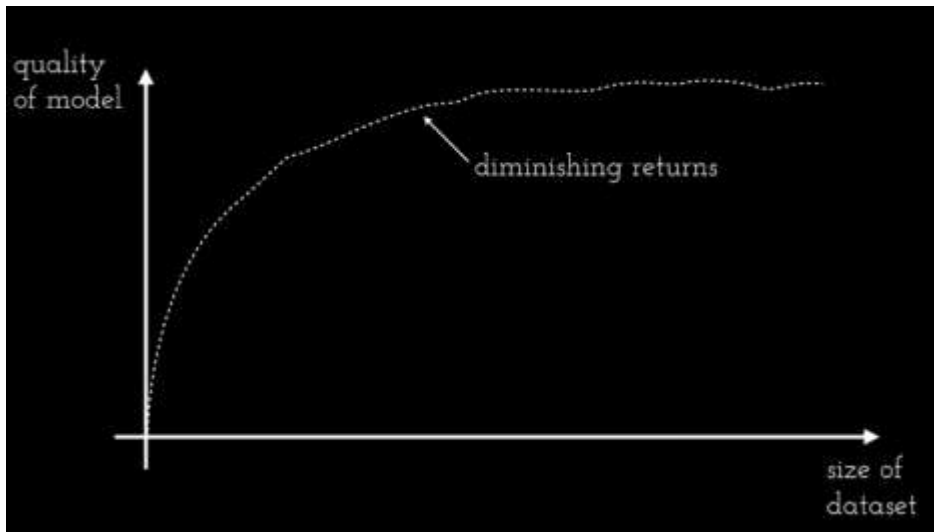
		Predicted									
		1	2	3	4	5	6	7	8	9	0
Actual	1	10					3			3	
	2		39		3		12				
	3	2		39				8			
	4	8			39			4	1	45	
	5		7			92		1			
	6		6				33		3	3	8
	7			8		8		43			
	8		2		3				32		
	9	8			8	4	5			89	
	0		8	8			8		8		46



Data Collection

More data beats smarter algorithms

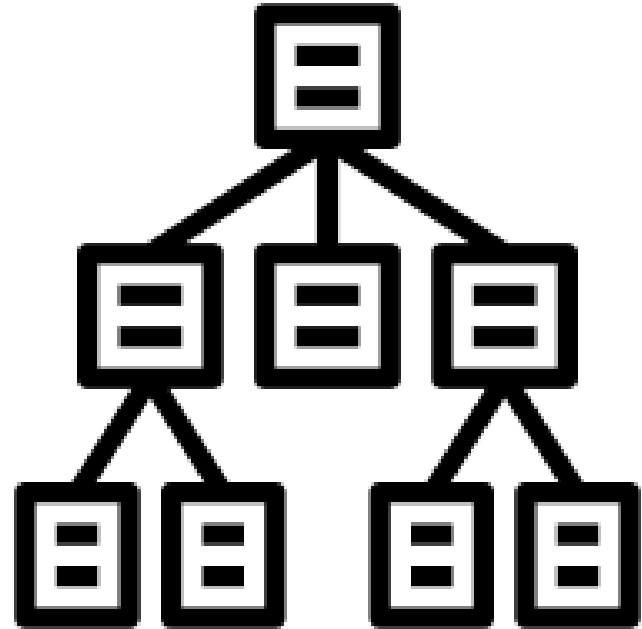
1. **But it is not practical**
2. Data is **expensive**. Money and time to collect labels
3. Big data might **be overkill**



Model Training

Based on different use cases

1. **Regression:** n dim-Polynomial?
2. **Classification:** Decision tree, logistic regression SVM
3. Each of the algorithm has multiple characteristics:
 - a. Susceptible to outliers
 - b. Explainability



Model Evaluation

Is it useful?

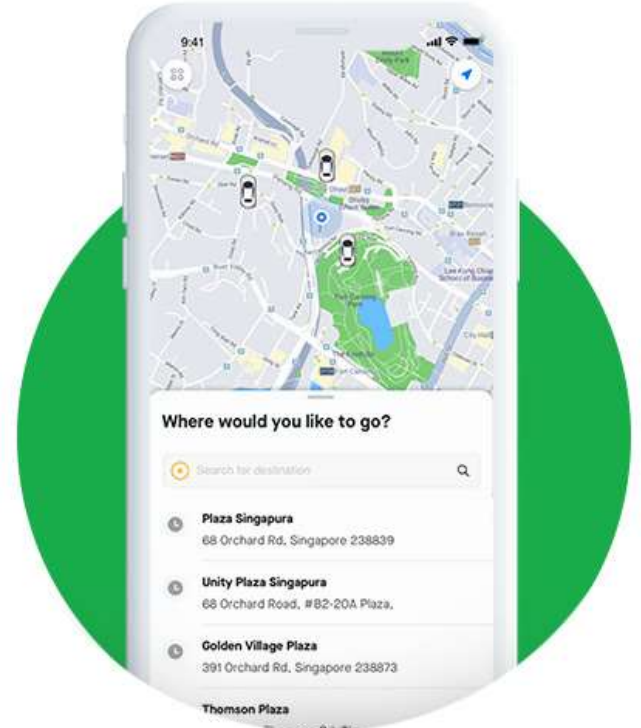
1. **Regression:** Root Mean Squared Error (RMSE)
2. **Classification:** Confusion metrics, AUC, Precision, Recall, F1
3. **Complexity, explainability, latency** (time and space)
4. **Eager/Lazy learners**

		Predicted									
		1	2	3	4	5	6	7	8	9	0
Actual	1	10					3			3	
	2		39		3		12				
	3	2		39				8			
	4	8			39			4	1	45	
	5		7			92		1			
	6		6				33		3	3	8
	7			8		8		43			
	8		2		3				32		
	9	8			8	4	5			89	
	0		8	8			8		8		46

ML Ops

Operating real ML for real Use Case

1. **Model Push**
2. **Model Validation**
3. **Monitoring/Anomaly Detection**



ML Tech Stack



Development + IDE

Proprietary + Confidential

The logo for Google Colab, featuring the word "colab" in a lowercase, rounded, orange font.

Data Studio



Language + Library

The pandas logo, featuring a stylized bar chart icon to the left of the word "pandas" in a bold, dark blue font.

Data + ML Ops



Google
BigQuery

The Apache Spark logo, featuring the word "SPARK" in a stylized font with a red star above the "A", and "APACHE" in smaller letters above it.

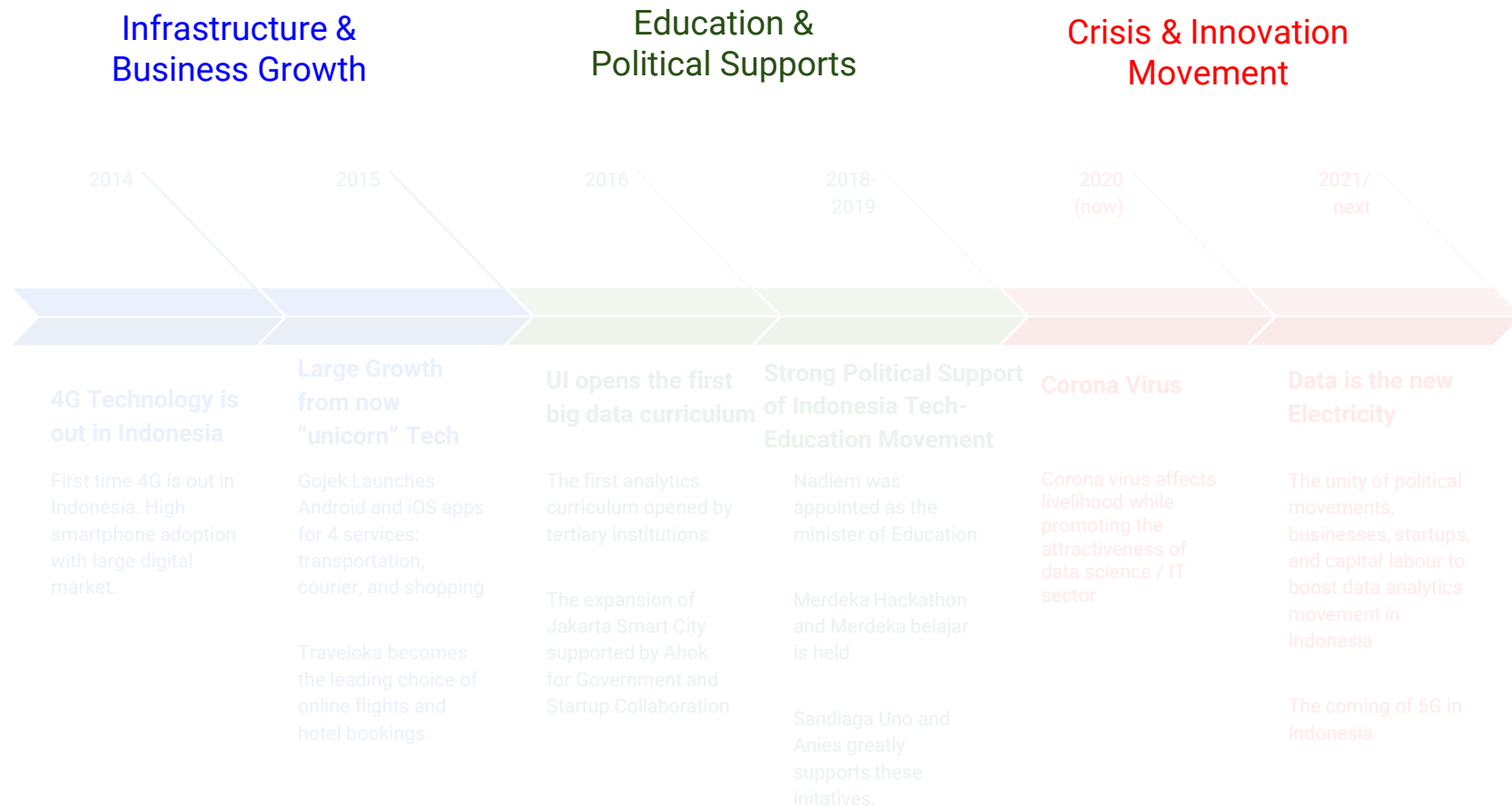
Analytics / ML Trend

How Analytics enter/menyurupi our lives?



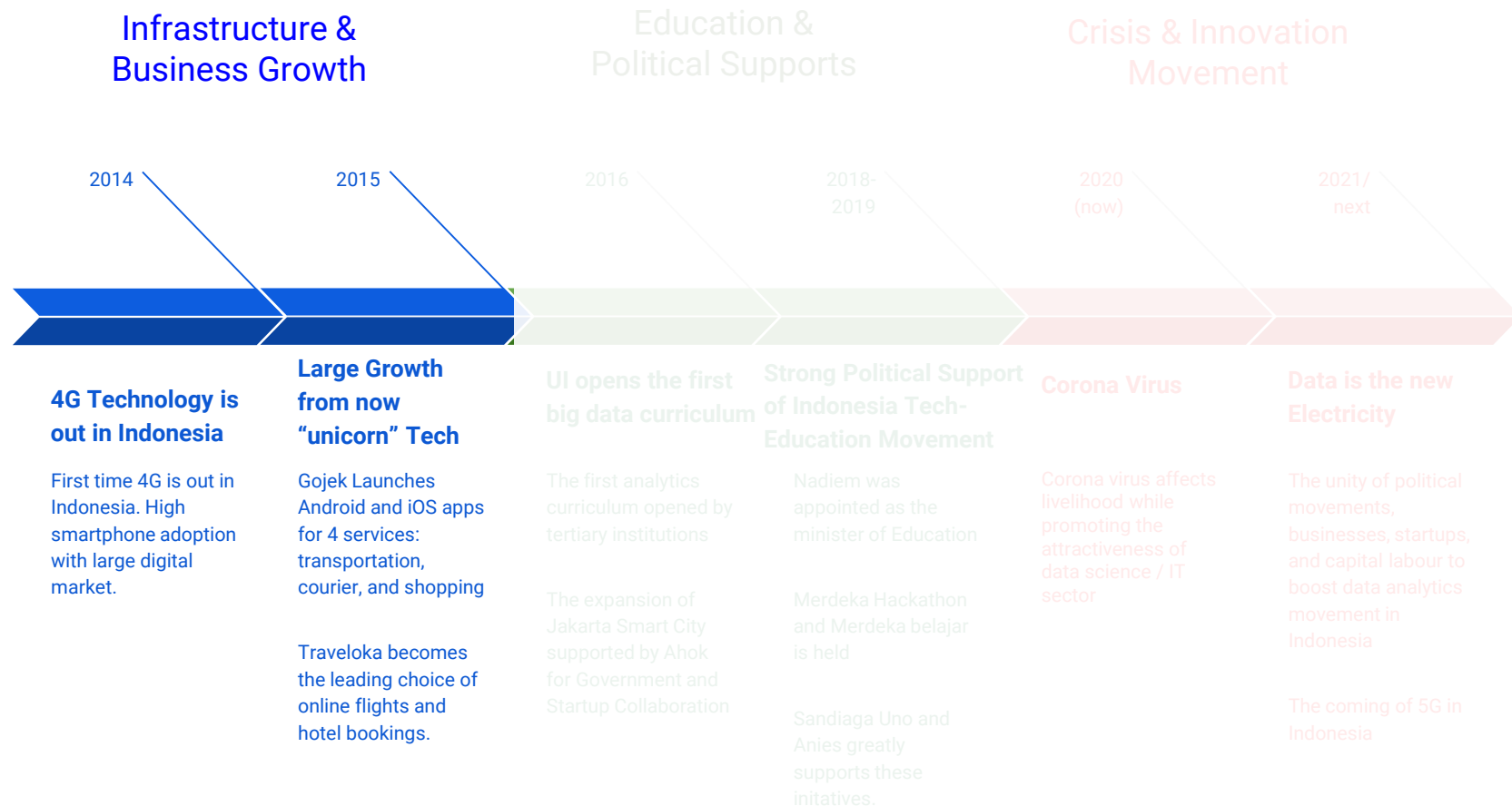
Analytics Development in Indonesia

Proprietary + Confidential



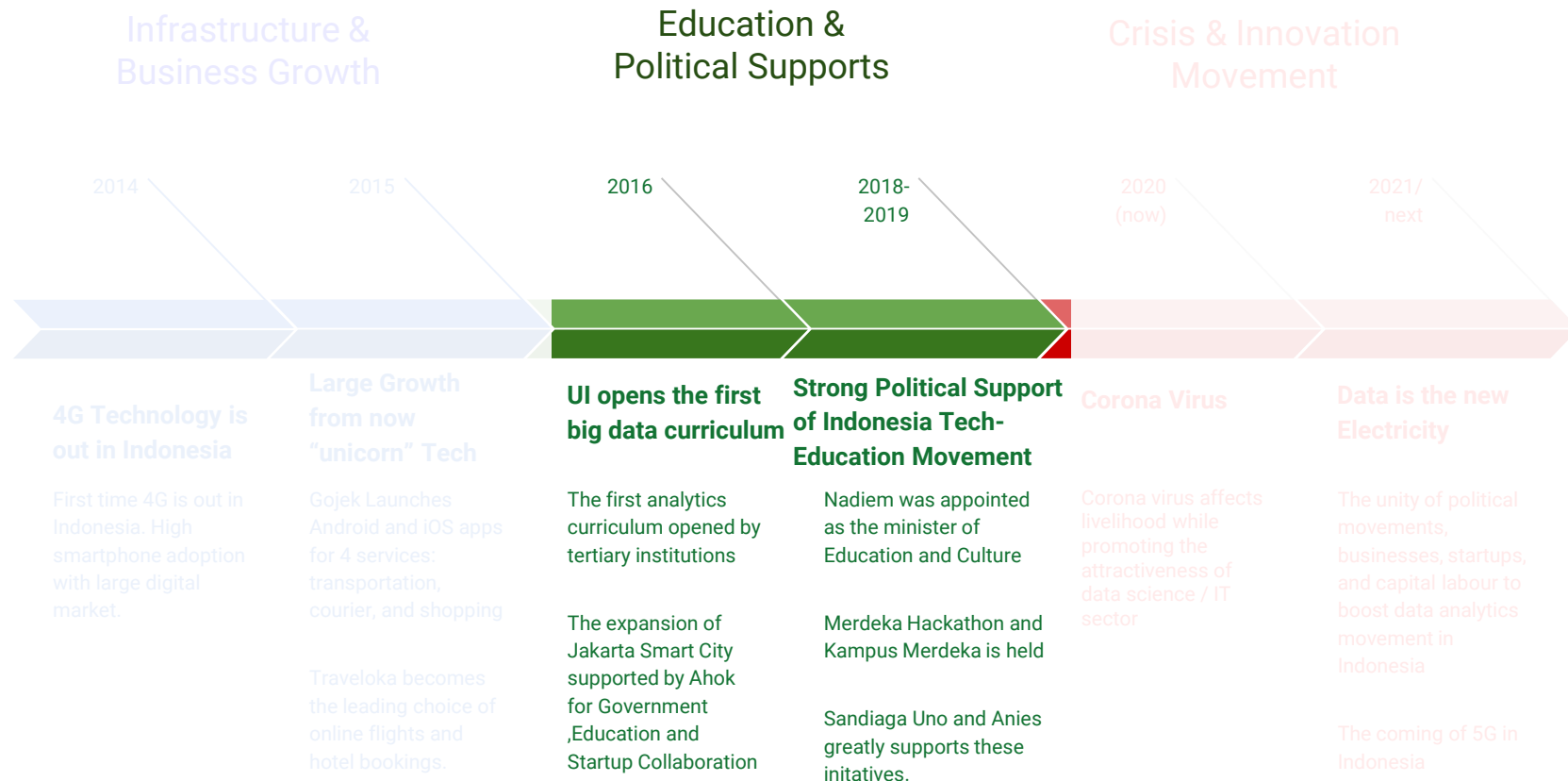
Analytics Development in Indonesia

Proprietary + Confidential



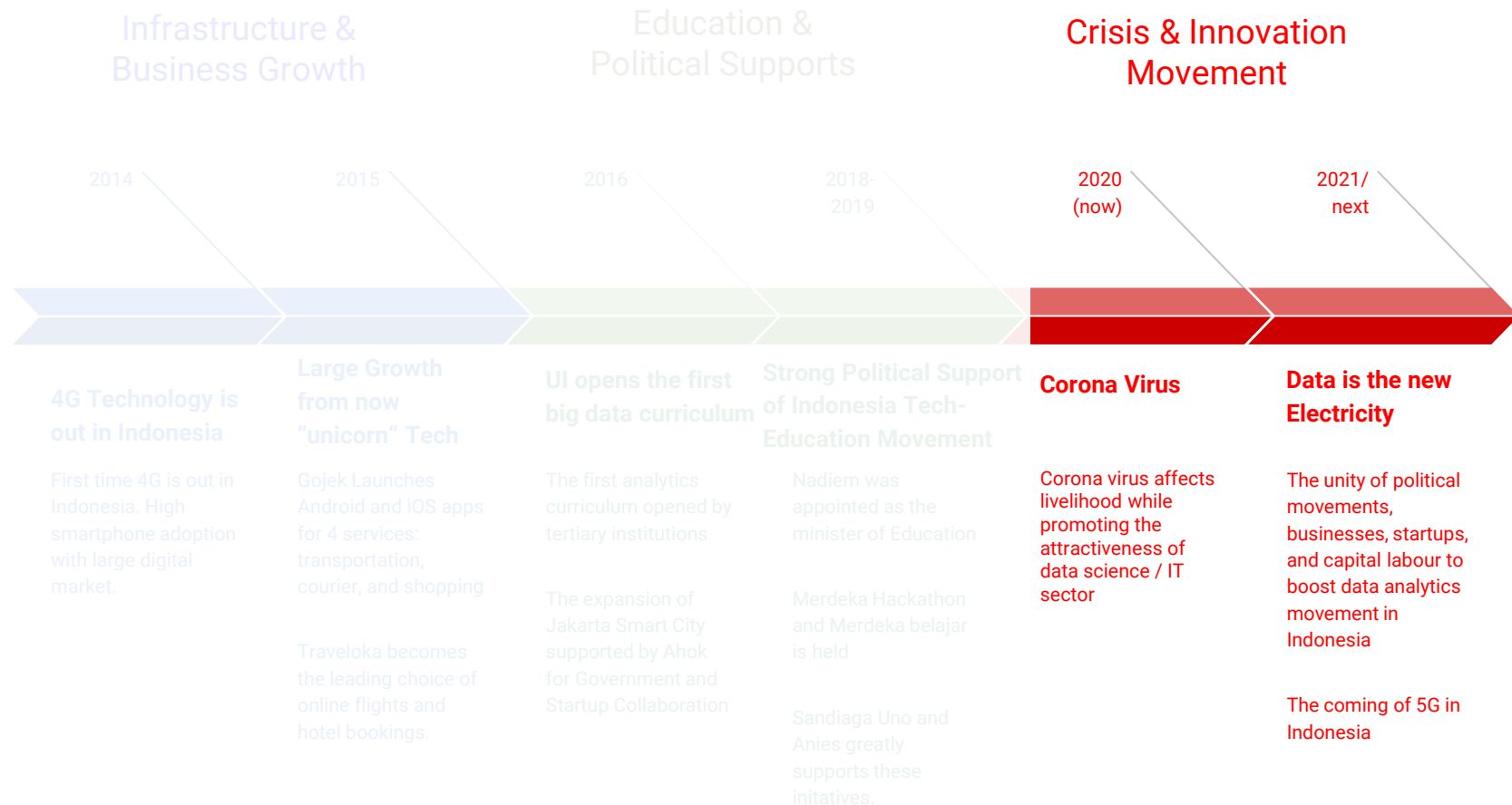
Analytics Development in Indonesia

Proprietary + Confidential



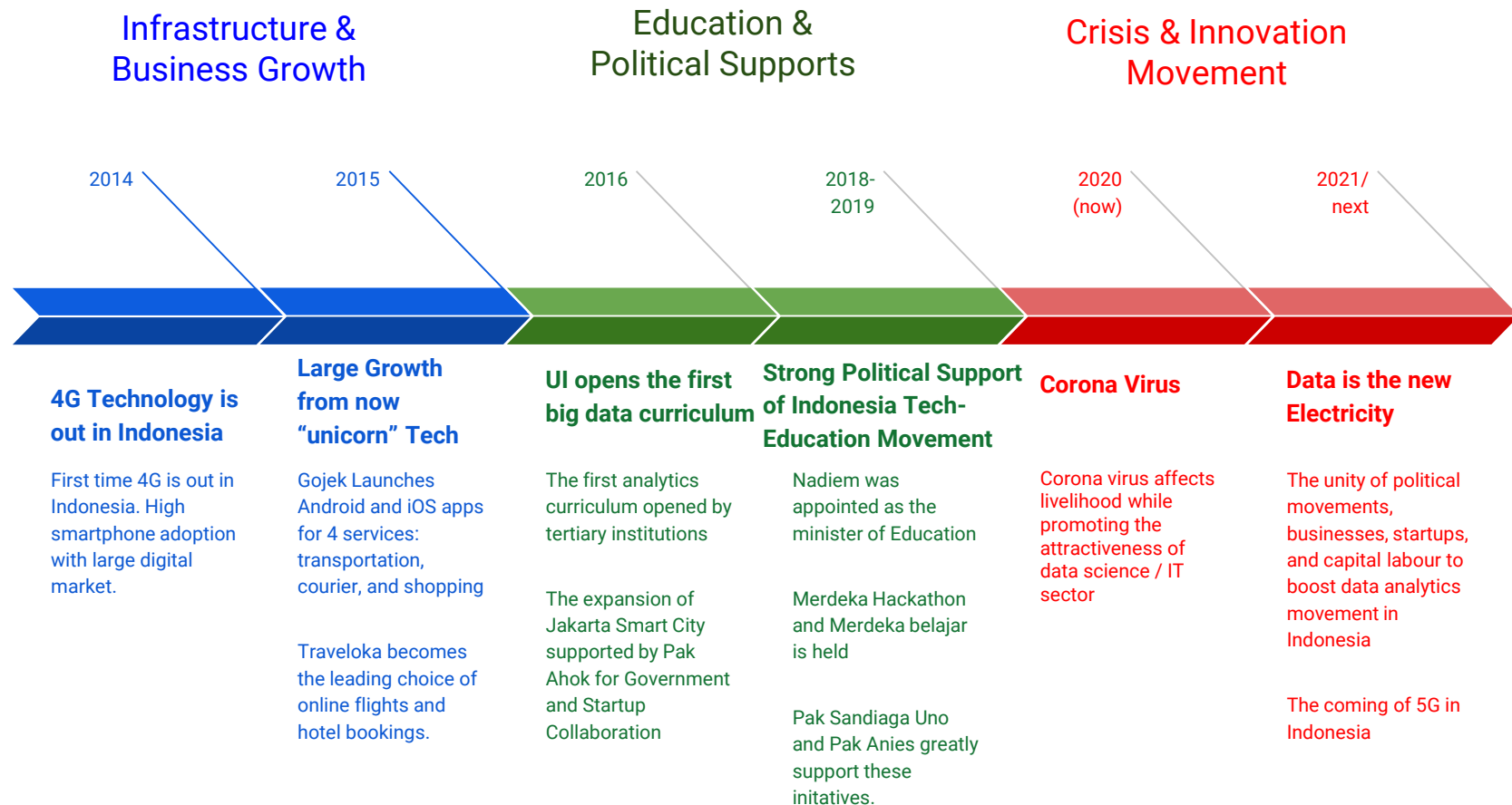
Analytics Development in Indonesia

Proprietary + Confidential



Analytics Development in Indonesia

Proprietary + Confidential



How can you excel in ML?

I'm super excited! What's next!



Data Analytics is hard... Here's how you can excel

What I have learnt over 3 years as Data Analysts/Engineers at Google, Visa, and Lazada *



Vincent Tatan

Oct 20 · 8 min read ★



Upskill yourself (source)

<https://bit.ly/2rBUP0Y>

Ace your Data Analytics Interviews

My experience Interview as a Data Analyst/ Scientist at Google, Visa and many more



Vincent Tatan

Sep 16 · 9 min read ★



<https://bit.ly/2Pc2axe>

Contribute more!

Data Science expands a lot, share your knowledge

1. **Read more:** Keep experimenting with your learning styles (kinesthetic, auditory, visual)
2. **Write more:** Write articles and share them!
3. **Speak more:** Teach your fellow peers or any conferences out there!

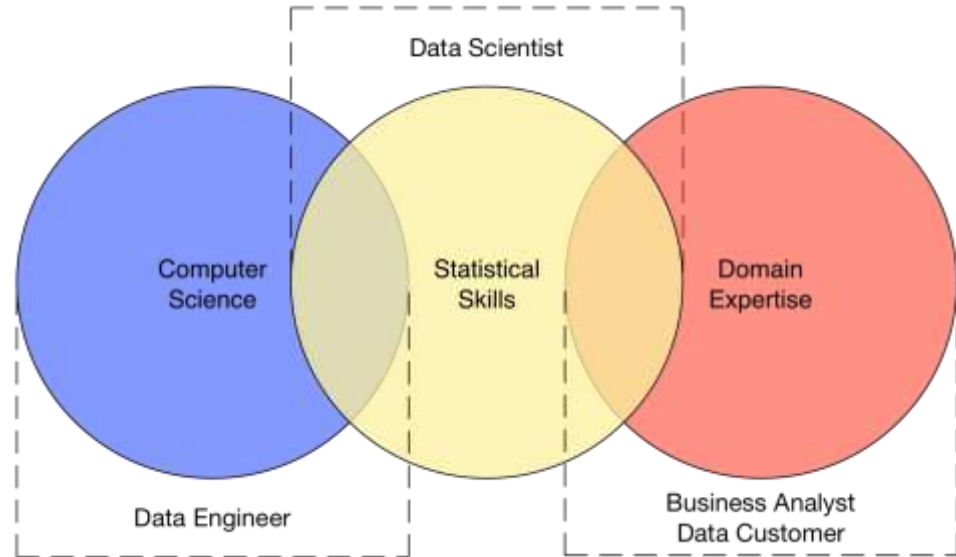
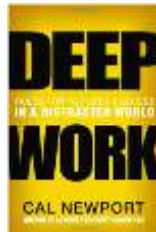


Generalize and Specialize

Follow it to the T

1. **Strength:** Utilize your biggest strengths
2. **Communicate:** Communicate your strength and impacts more.
3. **Learn in T:** SQL & Python/R are the breadth, then domain knowledge is your depth

[Read Deep Work](#)



Smile!

Data Science is Fun

1. **Play:** Tough, so have fun.
2. **Hack:** Use Saturdays to learn with friends.
3. **Celebrate impacts:** Data science is about building impacts. Start small and celebrate!



This is the key to excel

Contribute, Prepare, and Smile



Questions?

