

# Uber

**How can Uber apply  
Machine Learning to  
improve their Customer  
Experience?**

## TABLE OF CONTENTS

TABLE OF TABLES	7
TABLE OF DIAGRAMS	7
1.0 An improved quantitative pickup quality metric for uber	8
1.1 Signals	8
1.2 Metric	9
2.0 How can operators improve the uber experience?	11
2.1 Operators?	11
2.2 Increasing Time Estimation Accuracy	11
2.3 Reducing chance of driver loops	11
2.4 Increasing the rider/contract time	1
2.5 Decreasing steps walked by riders	1
3.0 Pickup automation and experience improvement through artifical intelligence & machine learning <sup>2</sup>	
	2
3.1 The business problem	2
3.2 Viability of Machine Learning	3
3.3 Data Gathering	4
3.4 Data Pre-processing	6
3.5 Machine Learning Algorithms	7
3.6 Features & Model improvement	7
3.7 User Experience	8
4.0 Bibliography	9

## TABLE OF TABLES

Table 1: Signal breakdown	8
Table 2: The benefits of a data-lake	4
Table 3: The Hudi datatypes to store the various features that could formulate the ML mode	5

## TABLE OF DIAGRAMS

Diagram 1: Visualisation of the metric's effect on the contextualise pain point	9
Diagram 2: Visualisation of the metric's effect on the depict pain point	10
Diagram 3: Proposed initial rider's screen while stuck in traffic	1
Diagram 4: Proposed new driver's screen	2
Diagram 5: Proposed rider's screen informing them of change in driver	2
Diagram 6: Proposed notification to rider	1
Diagram 9: The benefits to Uber of Machine Learning	2
Diagram 7: Typical Business Intelligence Architecture (Source: Chaudhuri, 2011)	4
Diagram 8: Uber's Business Intelligence Architecture (Source: Govindarajan, 2023)	5

# 1.0 An improved quantitative pickup quality metric for uber

## 1.1 Signals

To develop an accurate pickup quality metric, each feature's type must be conveyed:

Type of Signal	Potential features
Active	<ul style="list-style-type: none"><li>▪ User Generated Feedback</li><li>▪ Support tickets</li><li>▪ UX research</li><li>▪ Cancellation rate</li><li>▪ Rider/driver contact rate</li><li>▪ Pickup Location Error</li></ul>
Passive	<ul style="list-style-type: none"><li>▪ Sensor Inference</li><li>▪ Location scoring during uncertainty</li><li>▪ Driver loops around rendezvous</li><li>▪ Driver ETA jumps within last 300m</li><li>▪ ETA vs ATA</li><li>▪ Number of steps walked by riders</li><li>▪ Number of heading changes for driver</li><li>▪ Number of heading changes for rider</li><li>▪ Driver idle time</li><li>▪ Number of taps until request</li><li>▪ Time spent in pin edit</li></ul>
Third-Party	<ul style="list-style-type: none"><li>▪ Real-time congestion</li><li>▪ Hot spots</li><li>▪ X-axis, z-axis confidence swings</li></ul>

Table 1: Signal breakdown

Uber's business vision is to (Sawhney et al., 2019):

*Build an effortless pickup experience for everyone, everywhere, everytime*

Based on this vision the most important features are:

- ETA vs ATA
- Driver loops around rendezvous
- Rider/driver contact rate
- Number of steps walked by riders

The weighting shifts primarily towards passive signals, as these are guaranteed to be discrete values – active ones may require interpretation and translation (Qualaroo, 2023). In addition, human error or lack of effort could diminish active signal quality (Cloke, 2023).

Some features are subtly conveyed via the selected ones through reasoning, for instance:

ETA vs ATA being high implies either a PLE occurred or; a support ticket was filed

In another example: the number of steps walked being high could reveal a PLE.

## 1.2 Metric

Pickup quality = minimise (30% \* ETA vs ATA)  
+ minimise (20% \* Driver loops around rendezvous)  
+ minimise (25% \* number of steps walked by riders)  
+ **maximise** (15% rider/driver contact rate)  
+ minimise (10% \* sum(all third-party signals))

This selection can be corroborated utilising the journey-level pain points from the case study (Gibbons, 2021; Sawhney et al., 2019):

The *contextualise* pain point reveals that a low-confidence location resulted in 50 steps being walked by a rider with foot pain. If all else was positive, then reduction of this feature - via a better location-confidence; would result in a far better experience.

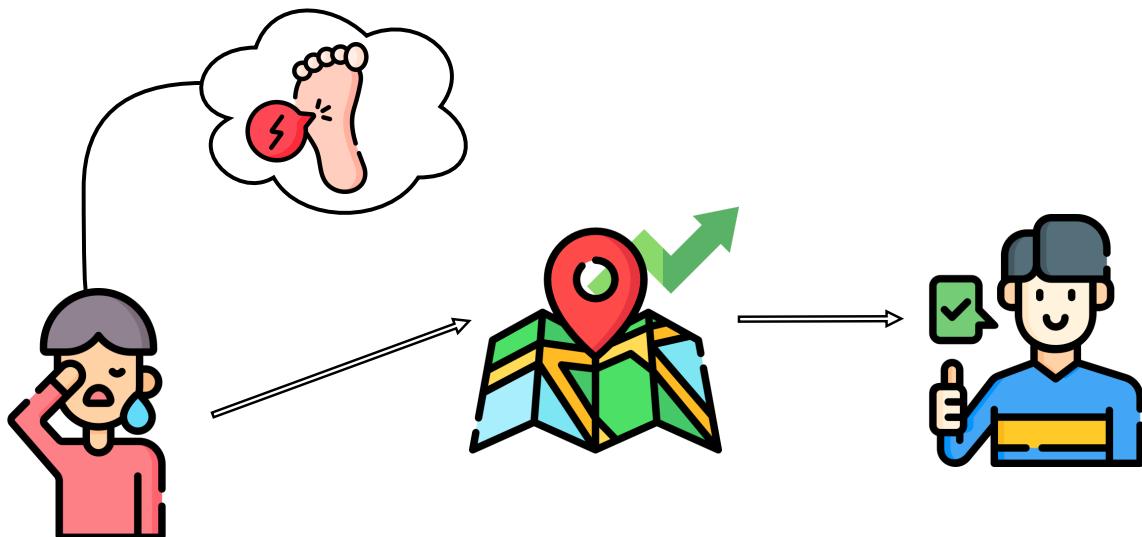
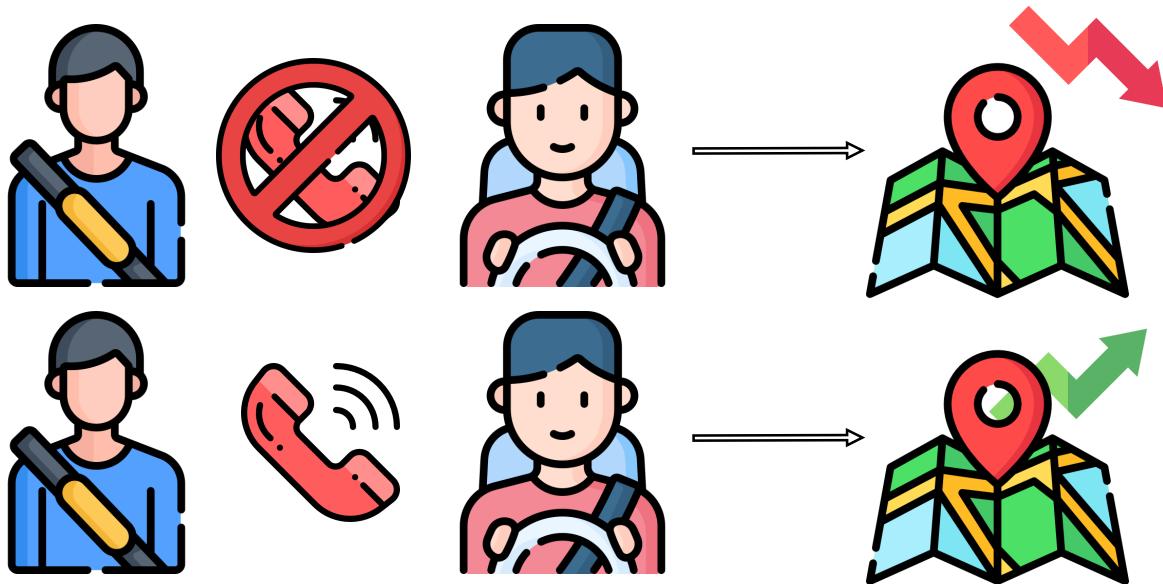


Diagram 1: Visualisation of the metric's effect on the contextualise pain point

Likewise, the *depicts* scenario showcases a good pickup experience. However, to elevate it: maximising the rider/driver contact would have facilitated communication of the exact rider location.



*Diagram 2: Visualisation of the metric's effect on the depict pain point*

The metric should accurately represent Uber's *north star metric*: the number of successful rides completed per day (Ellis, 2017).

Strengths:	Weaknesses:
<ul style="list-style-type: none"> <li>▪ Recognised brand</li> <li>▪ Little competition</li> <li>▪ High standard</li> <li>▪ Rating system</li> <li>▪ Low price</li> <li>▪ Proxy drivers</li> </ul>	<ul style="list-style-type: none"> <li>▪ No personal relationships</li> <li>▪ Unpredictable</li> <li>▪ No fool-proof driver vetting</li> <li>▪ Data protection risk</li> </ul>
Opportunities:	Threats:
<ul style="list-style-type: none"> <li>▪ Dissatisfaction with taxis</li> <li>▪ Open to new countries</li> <li>▪ Transition to electric cars</li> </ul>	<ul style="list-style-type: none"> <li>▪ Low wages</li> <li>▪ Could allow new company into market</li> <li>▪ Self-driving cars</li> <li>▪ Passenger safety legality</li> </ul>

*Table 2: SWOT analysis of Uber*

It should be noted that if the feature range was not restricted, numerous could be added:

- Congested routes taken could be minimised
- Driver denial of service due to intoxication could be recorded
- Driver to rider score-matching could be employed
- Use of voice-commands should be recorded – hotspots could see a surge in usage
- Phone service providers could be partnered with to improve low-confidence locations when GPS is weaker

## 2.0 How can operators improve the uber experience?

### 2.1 Operators?

For clarity of argument, it is inferred that operator describes any user that engages with the active uber experience.

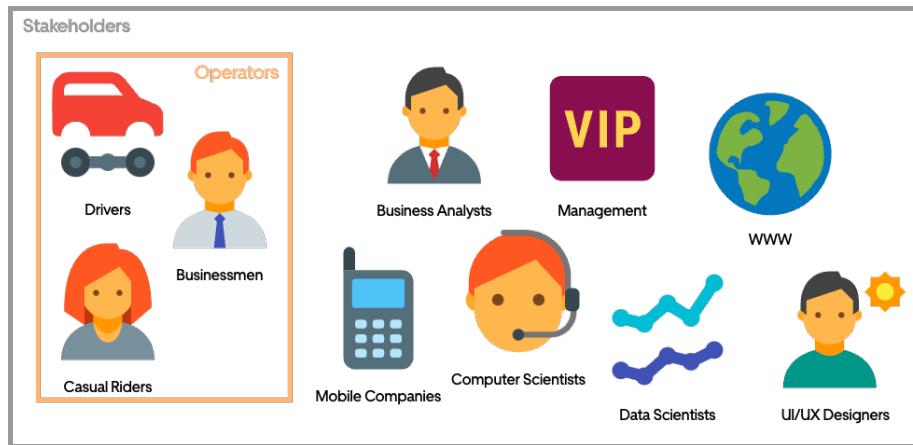


Diagram 3: Scope of "operators" for future reference

### 2.2 Increasing Time Estimation Accuracy

Uber's estimated arrival-times are generated utilising predictive analytics, with accepted error margins (Hu et al., 2022a). Uber utilises terrain maps and the Global Navigation Satellite System [GNSS] (Iland, 2018). Certain functions are employed to minimise prediction errors (i.e. *shadow mapping* allows for accuracy in low-confidence areas). While this is beneficial; the model currently requires downtime to test, limiting real-time accuracy (Hu et al., 2022b). A method should allow for *continuous, incremental*/training instead (Hu et al., 2022b).

### 2.3 Reducing chance of driver loops

Drivers looping could be reduced by increased rider/driver contact time, and extreme options:

If a driver takes too long to arrive at their destination, and has been static for a while, a notification could allow them to re-offer the job to nearby ubers. The driver could receive a percentage compensation, or still complete the job if they begin to move again – however it a car that is on a less congested route could complete the job in a more timely fashion: keeping uber customer-centric.

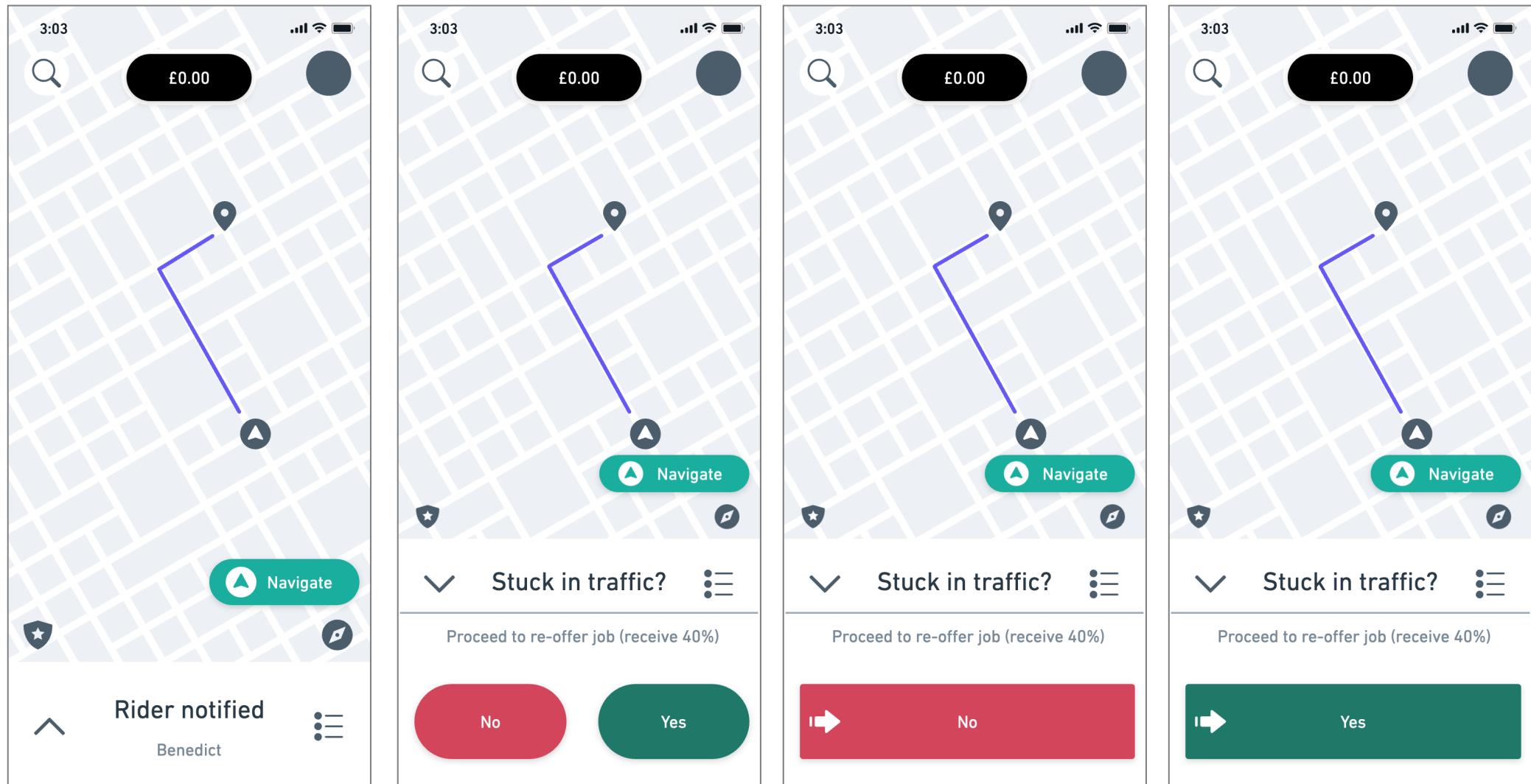
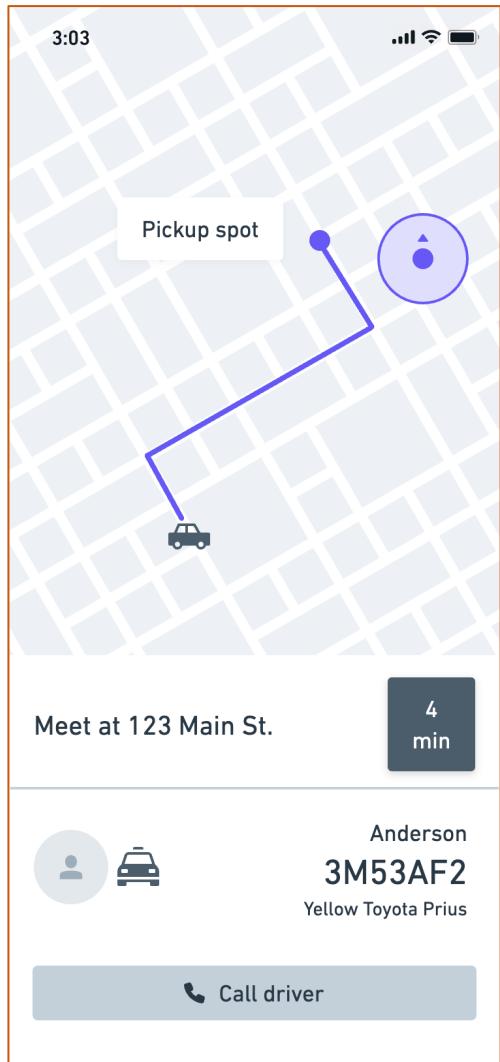


Diagram 4: Proposed initial rider's screen while stuck in traffic

Count: 3083

20187281

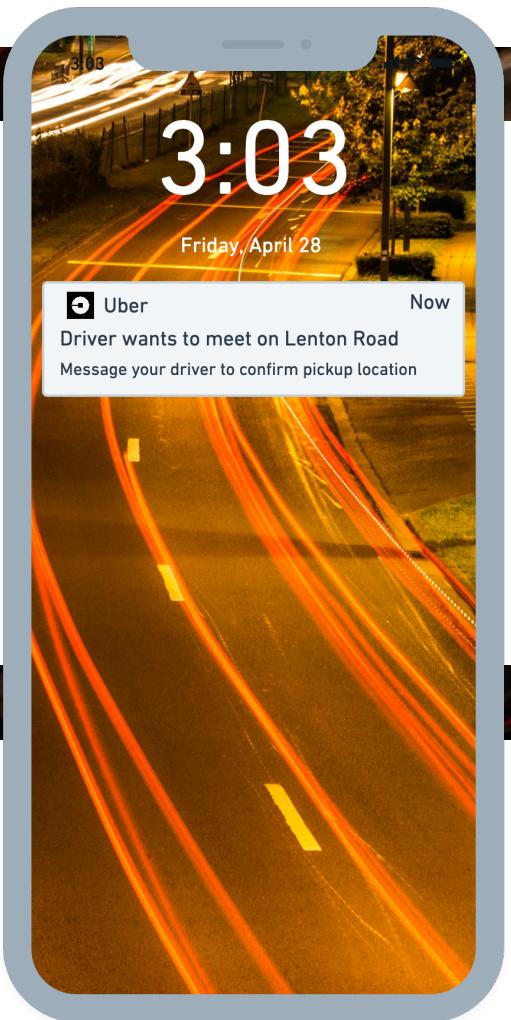


In addition, a slight delay to movement animation on the rider side application is recommended. The map view is prone to car icon jumps in areas when *shadow mapping* is applied – instead of just GNSS (Shekhar & Sankaravadiel, 2019). Adding a subtle delay would gift more time to allow location smoothing (Shekhar & Sankaravadiel, 2019).

## 2.4 Increasing the rider/contract time

To increase contact time between operators: drivers must either text or call riders upon arrival **at least** once. If not, it could negatively affect their score. Additionally, a notification should tell the rider to message their driver, not just meet them.

Diagram 7. Proposed notification to rider



## 2.5 Decreasing steps walked by riders

Currently, map-matching calculates accurate spots for drivers (Shekhar & Sankaravadiel, 2019). Other methods are more accurate but computationally intense (Shekhar & Sankaravadiel, 2019). An example is *dead reckoning* which allows for smoothing between poor GPS points utilising mathematical equations (Shekhar & Sankaravadiel, 2019). To conglomerate these results: Beacon conveys information to differing phones by deciding on formulae in accordance with the device in use (Shekhar & Sankaravadiel, 2019).

Some devices do not convey certain parameters for shadow matching – as they have rules governing information sharing. For instance, iPhones do not automatically send information, thus other, costly methods must be expended [i.e. Bluetooth connections] (Shekhar & Sankaravadiel, 2019).

iPhones thus need to be better leveraged technologically to extract their GPS information, to allow for more time to correct errors and display information to the user.

### 3.0 Pickup automation and experience improvement through artifical intelligence & machine learning

#### 3.1 The business problem

How does one provide *an effortless pickup experience for everyone, everywhere, everytime*

The basic need offered by Uber is effortless transit (Cooper, 2014; Sawhney et al., 2019): This benefits the customer as well as the riders - who gain an income avenue; while external stakeholders gain profits.

The desired objective is provision of the perfect pickup (MBEC, 2022; Sawhney et al., 2019). Steps towards combatting the business problem have already been initiated, with a 2016 paradigm-shift towards destination-first customer autonomy (Sawhney et al., 2019). In addition, two ML platforms have been employed before: *Horovod* and *Michelangelo* – but are open to improvement (Shao & Islam, 2021a). Uber actively embraces Artificial Intelligence, so the goals of this solution are in alignment with the company's philosophy. In conjunction, the benefits of such decision support will be monetary in value. (Dean, 2014):

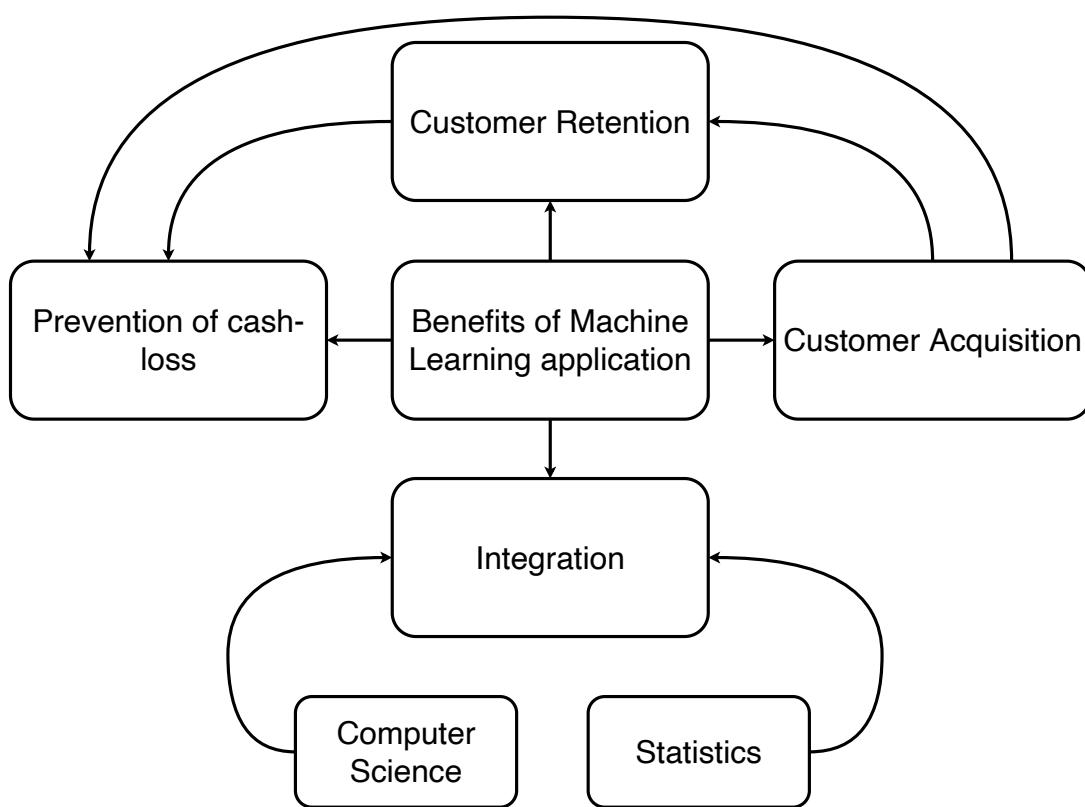


Diagram 8: The benefits to Uber of Machine Learning

Uber's true problem is the need to enhance methods for relationship management between stakeholders. The principle of choice henceforth is a model that can provide said enhancement consistently (Turban et al., 2015). In solving the surface and deep-level business problems;

Uber could adopt the Agile methodology for *continuous improvement* upon the ML model that is to be described ([Project Management Institute, 2017; Harirajan, 2022](#)).

## 3.2 Viability of Machine Learning

Machine learning, which has already been employed, is only of use – in this scenario – to Uber if other methods of prediction fail (Dean, 2014). Simple mathematical equations are not suitable for the number of factors that Uber stores. This data is largely unstructured, stored in a format that is difficult in understanding for humans, crafting a perfect opportunity for ML to be employed (Lee & Obermeyer, 2017).

Currently, Uber utilises a metric to model the quality of their data storage – CAP (Zhong & Cheng, 2021):

$$\text{Cost-Efficiency} \times \text{Accuracy} \times \text{Performance} = \text{Constant}$$

The underperformer in the model currently is performance, with the cost of collecting data being minimised already by a number of size reduction techniques such as – data storage on *the Apache® Hadoop® File System (HDFS) space*, in *Apache Hive™ tables* (Shao & Islam 2021b).

Machine learning is necessary due to the sheer amount of data that Uber receives. It has 134,000 HDDs allocated to information storage, with the average size being 4TB (Shao & Islam, 2021c).

ML is therefore a viable option: there is a well-defined problem to be solved by such algorithms; there is an abundance of continuously updated data; said data is relevant to the problem and cleanable; the results of data analysis would directly impact the problem and; it is clear that the volume of data prevents problem solution through traditional means (Dean, 2014).

A structured decision problem can thus be presented:

How can profit be maximised, when profit directly correlates to pickup quality

&

Pickup quality = minimise (30% \* ETA vs ATA)  
+ minimise (20% \* Driver loops around rendezvous)  
+ minimise (25% \* number of steps walked by riders)  
+ maximise (15% rider/driver contact rate)  
+ minimise (10% \* sum(all third-party signals))

It must be acknowledged, that the use of passive signals is to mitigate the subjectivity of some active signals .

### 3.3 Data Gathering

Uber's data storage initially utilised only Online transaction processing databases [OLTP] to handle their data requirement (Reza, 2018). This provided fast latency for TBs of data. In 2016, this was replaced with a data-warehouse. Difficulty in its use stemmed from a lack of rule enforcement on JSON entries from data sources (Reza, 2018).

Following this, Hadoop adoption facilitated data-lake creation: to store standardised Parquet data (Reza, 2018). Since, Hudi has been added to abstract and allow querying based on time metrics (Reza, 2018). This data-lake is a colossal information store, combining features of both operational databases and data-warehouses via a Kimball approach (Reza, 2018; Sharma, 2023a; Govindarajan, 2023):

Operational database	Data-warehouse	Data-lake
<ul style="list-style-type: none"> <li>▪ Concerns current data</li> <li>▪ Utilised by database professionals</li> <li>▪ Useful for business maintenance</li> </ul>	<ul style="list-style-type: none"> <li>▪ Concerns historical data</li> <li>▪ Utilised by managers &amp; analysts</li> <li>▪ Useful for data analysis</li> </ul>	<ul style="list-style-type: none"> <li>▪ Concerns both historical and current data</li> <li>▪ Utilised by a multitude of professionals</li> <li>▪ Useful for data analysis, business management, and data visualisation</li> </ul>

Table 3: The benefits of a data-lake

Hudi stores data in adherence with the ACID transaction methodology to ensure data integrity (Agarwal, 2020; IBM, 2023). Schemaless then stores Uber trip data (Thomsen, 2016). Of note, is the lack of emphasis on the transformation portion of the ETL framework –preferring raw data to manipulate (Reza, 2018).

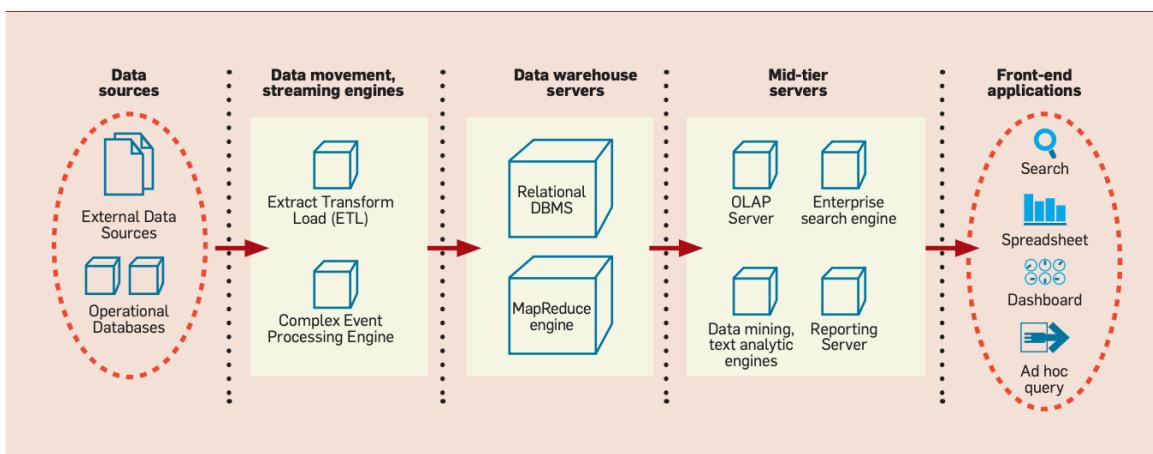


Diagram 9: Typical Business Intelligence Architecture (Source: Chaudhuri, 2011)

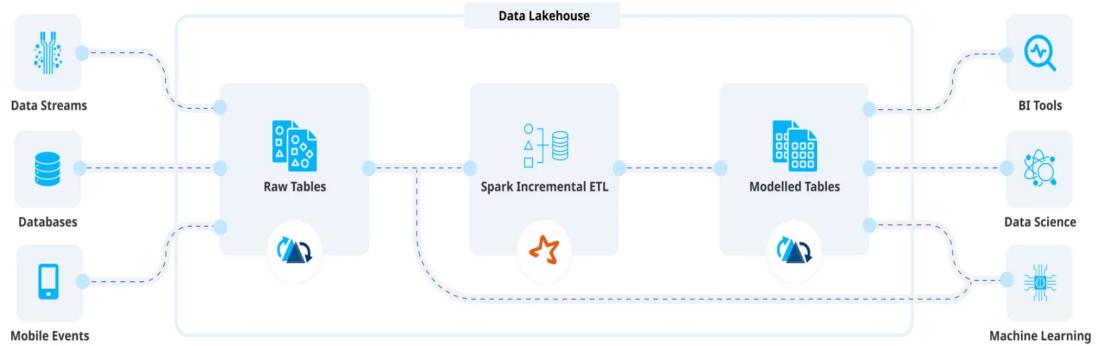


Diagram 10: Uber's Business Intelligence Architecture (Source: Govindarajan, 2023)

Through the data-lake, Uber collects – at the very least – this data which can be employed in the machine learning model (Sawhney, 2019):

Features	Datatype
<ul style="list-style-type: none"> <li>▪ User Generated Feedback</li> <li>▪ Support tickets</li> <li>▪ UX research</li> <li>▪ Cancellation rate</li> <li>▪ Rider/driver contact rate</li> <li>▪ Pickup Location Error</li> </ul>	<ul style="list-style-type: none"> <li>⇒ String</li> <li>⇒ String</li> <li>⇒ String</li> <li>⇒ Float</li> <li>⇒ Float</li> <li>⇒ Float</li> </ul>
<ul style="list-style-type: none"> <li>▪ Sensor Inference</li> <li>▪ Location scoring during uncertainty</li> <li>▪ Driver loops around rendezvous</li> <li>▪ Driver ETA jumps within last 300m</li> <li>▪ ETA vs ATA</li> <li>▪ Number of steps walked by riders</li> <li>▪ Number of heading changes for driver</li> <li>▪ Number of heading changes for rider</li> <li>▪ Driver idle time</li> <li>▪ Number of taps until request</li> <li>▪ Time spent in pin edit</li> </ul>	<ul style="list-style-type: none"> <li>⇒ String</li> <li>⇒ Float</li> <li>⇒ Int *</li> <li>⇒ Int *</li> <li>⇒ [Two different timestamps]</li> <li>⇒ Int *</li> <li>⇒ Int *</li> <li>⇒ Int *</li> <li>⇒ [Two different timestamps]</li> <li>⇒ Int *</li> <li>⇒ [two different timestamps]</li> </ul>
<ul style="list-style-type: none"> <li>▪ Real-time congestion</li> <li>▪ Hot spots</li> <li>▪ X-axis, z-axis confidence swings</li> </ul>	<ul style="list-style-type: none"> <li>⇒ Boolean + [Collection of geospatial values]</li> <li>⇒ Boolean + [Collection of geospatial values]</li> <li>⇒ [Multiple Floats]</li> </ul>

Table 4: The Hudi datatypes to store the various features that could formulate the ML mode

\*In another filesystem these would be stored as smallInt or byte – to reduce memory – but Hudi's specification allocates these to int

The data collected for the quality metric offsets any trade-offs that may exist due to its direct applicability to the business problem.

In conjunction, the inclusion of personal input is limited, due to the risk of bias affect predictions (Matthews, 2018).