

Bringing ML Online: Defining ML Features in Real-Time with Bytewax

Zander Matheson, Founder & CEO @  **bytewax**



Recommendations	Addictive TikTok
Ads	Annoying or not, at least less so with real-time relevance
Personalization	Mind reading in 2023, or at least it feels like it
Fraud Detection	Alerting us of our suspect train ticket purchases in foreign countries 😊

I'm Zander!

Working on Bytewax -> github.com/bytewax/bytewax

Proud human and dog dad

This photo is me trying to look cool 😏

You can find me in the Bytewax slack, or if you are in Santa Cruz, send me a LinkedIn and we can grab a coffee.



Today's Agenda

1. **What and why real-time ML?**
2. **Building a pipeline to analyze streaming data**
3. **Building a real-time feature pipeline**

- 1. What and why real-time ML?**
- 2. Building a pipeline to analyze streaming data**
- 3. Building a real-time feature pipeline**



Google	For every 100 ms of latency -> drop in search traffic of ~0.2%
Akamai	For every 100 ms of latency -> a 7% drop in conversion rates
Walmart	For every 1s of load time improvement -> a 2% increase in conversion rates
Staples	For every 1s of load time improvement -> a 10% increase in conversion rate



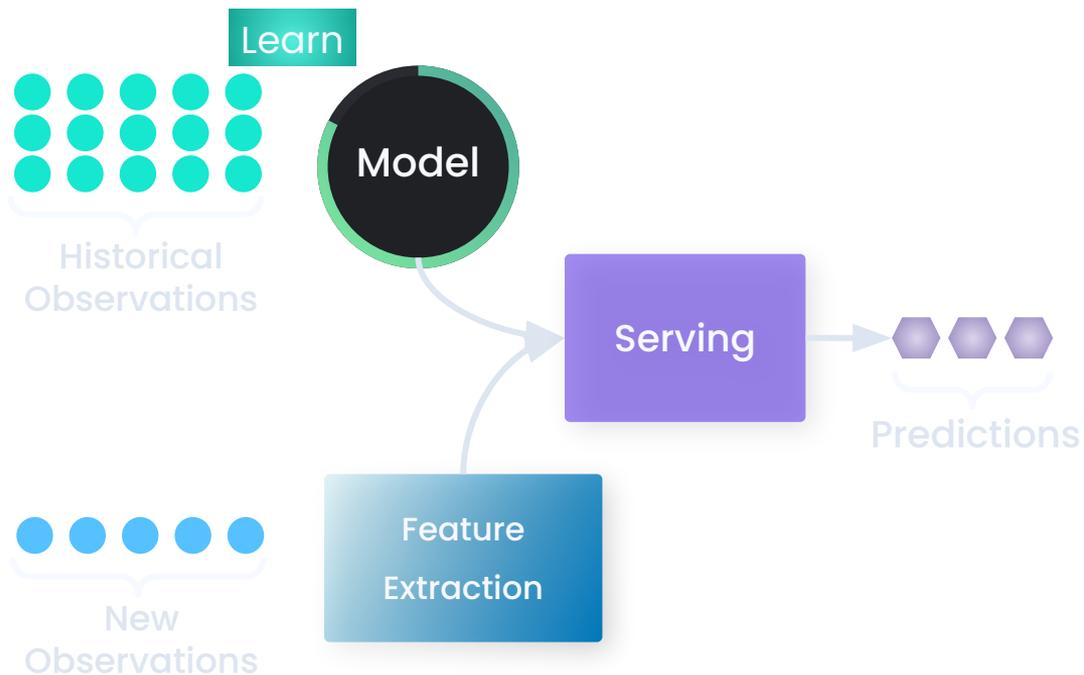
Real-Time Machine Learning

Real-Time ML can be classified into to types:

1. Offline training with online serving
2. Online training **and** serving

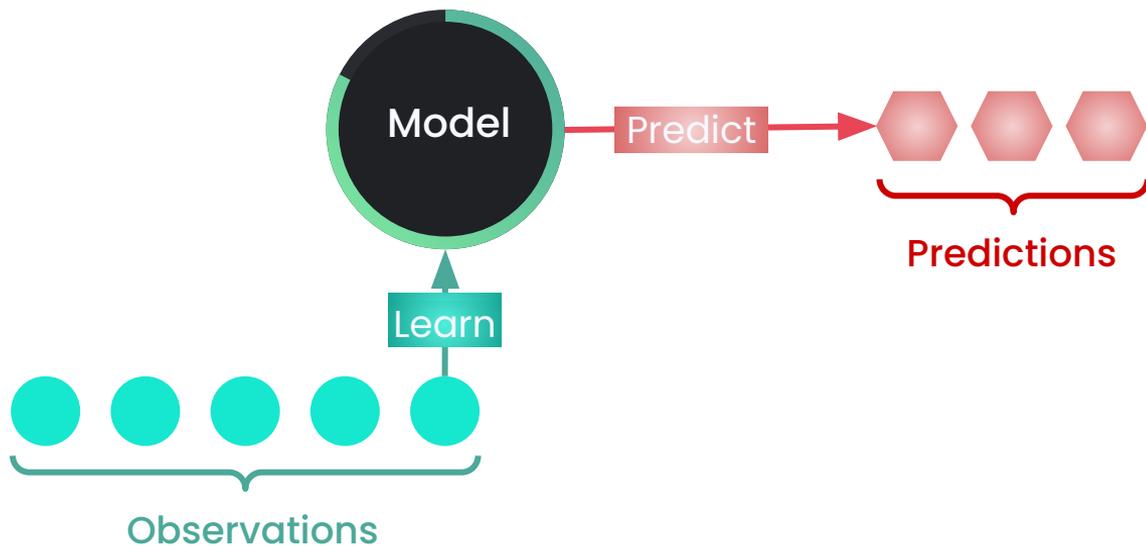


Offline Training, Online Serving



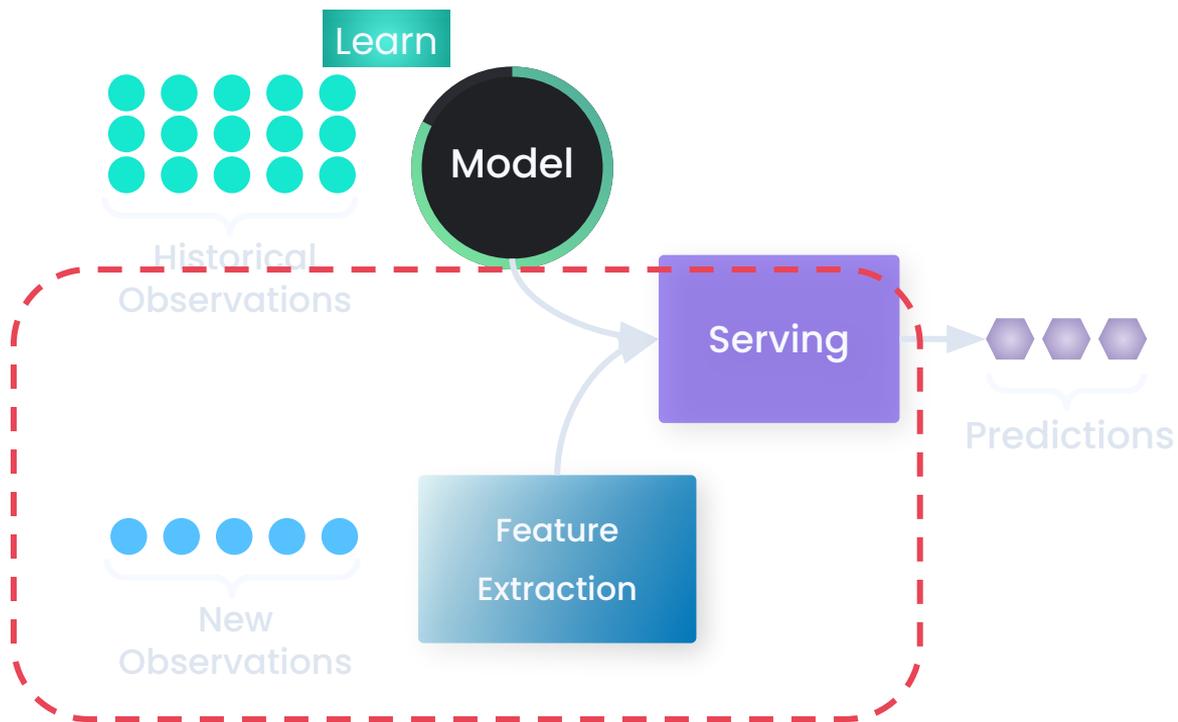


Online Training and Serving





Offline Training, Online Serving

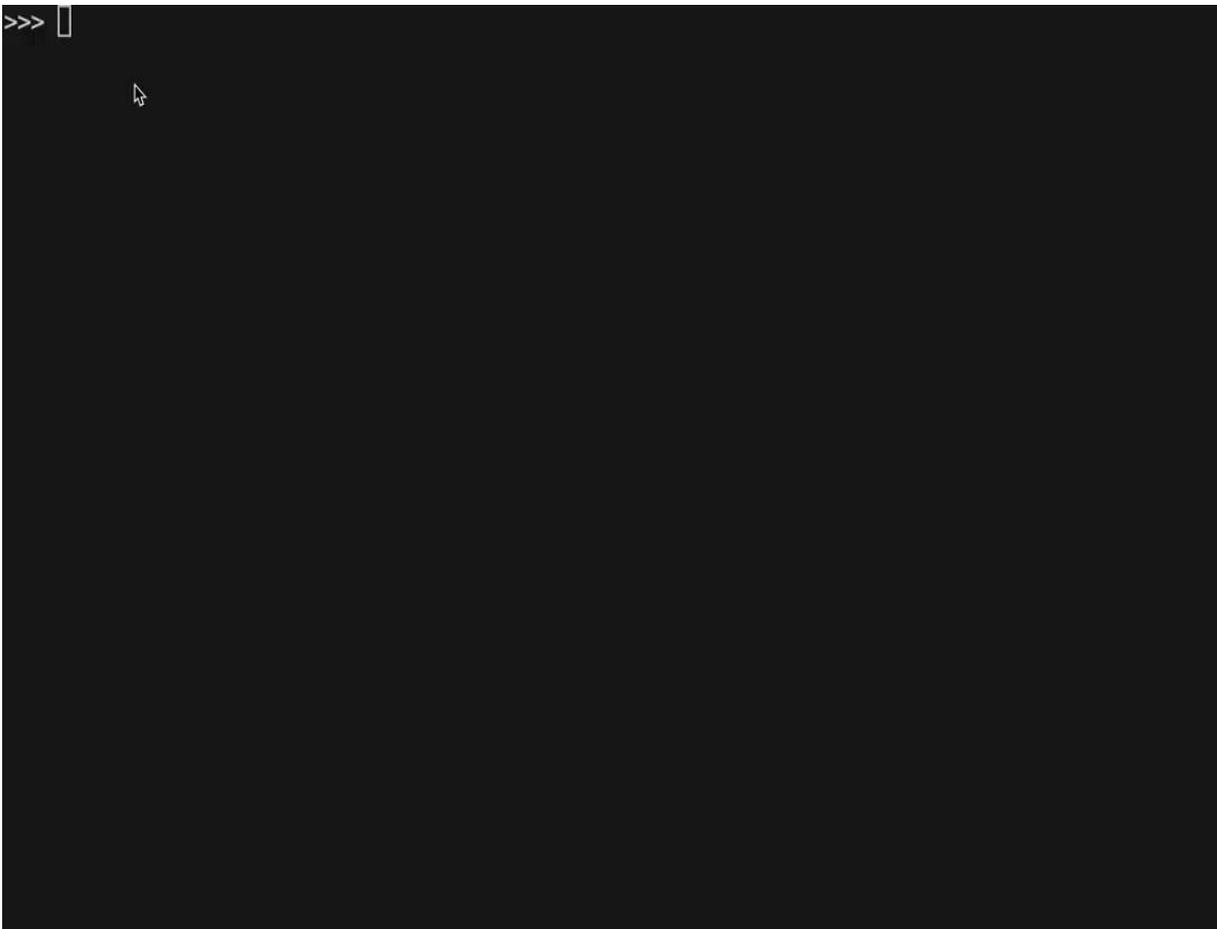


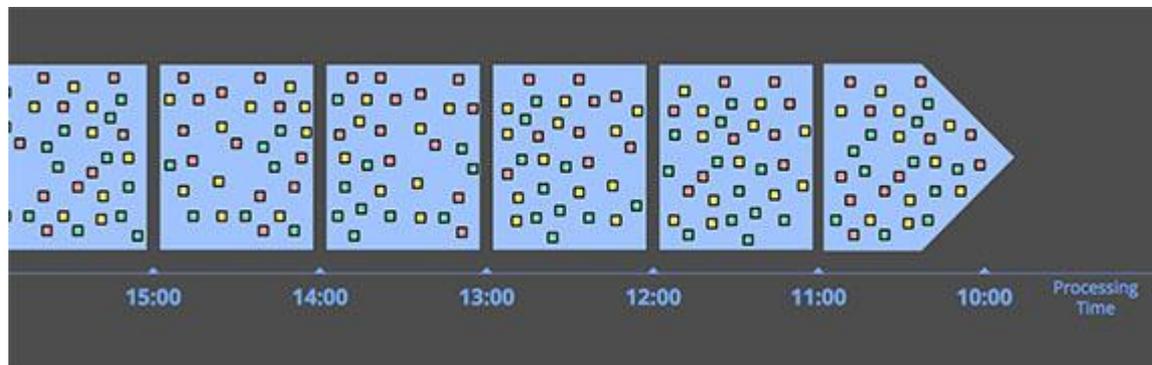
Outline:

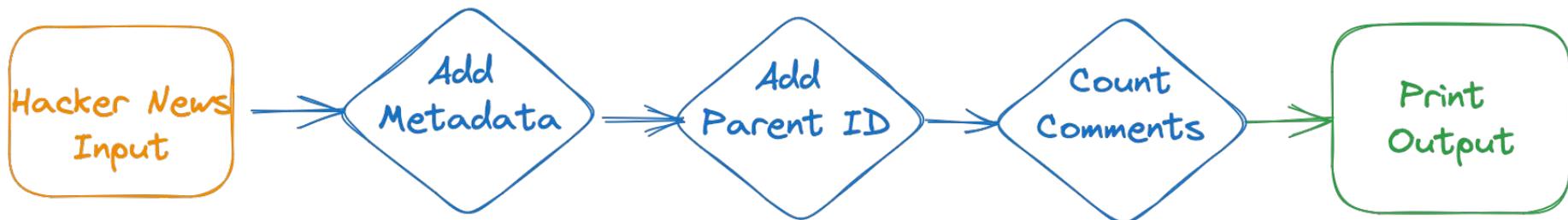
1. What and why real-time ML?
2. Building a pipeline to analyze streaming data
3. Building a real-time feature pipeline



Streaming Data







```

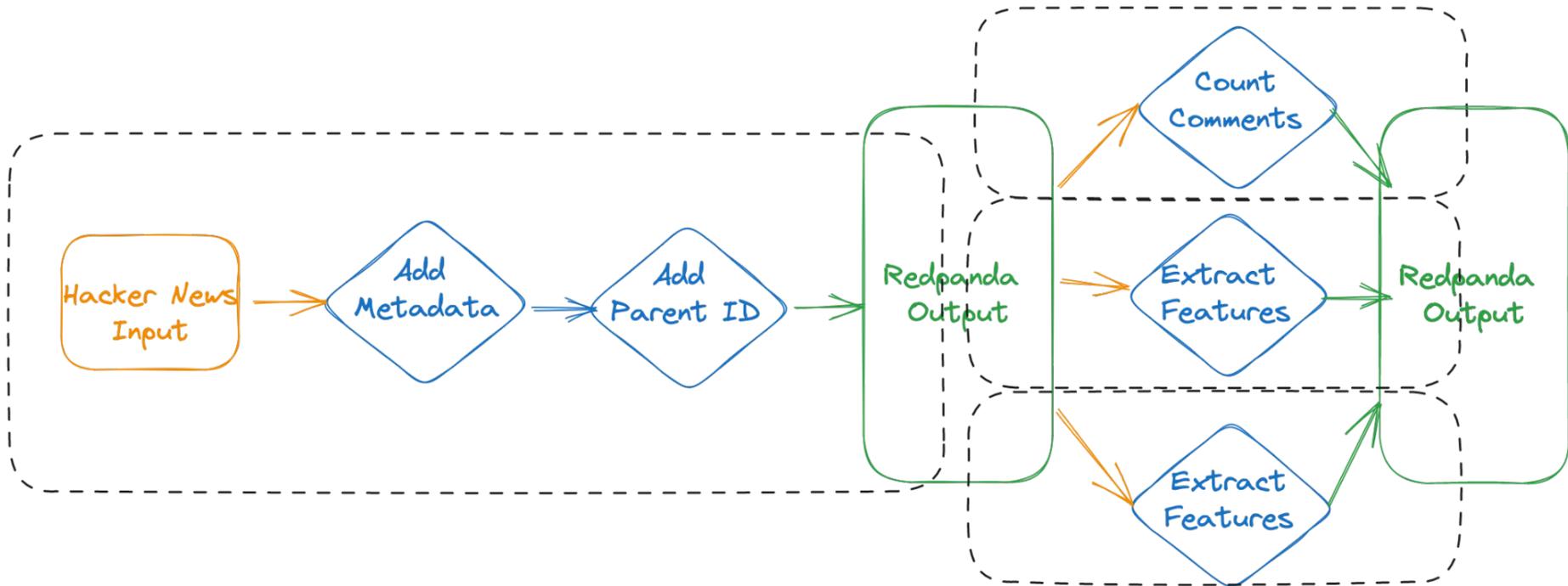
1 import requests
2 from datetime import datetime, timedelta
3 import time
4 from typing import Any, Optional
5
6 from bytewax.connectors.periodic import SimplePollingInput
7 from bytewax.connectors.stdio import StdOutput
8 from bytewax.dataflow import Dataflow
9
10 import logging
11
12 logging.basicConfig(level=logging.INFO)
13 logger = logging.getLogger(__name__)
14
15 class HNInput(SimplePollingInput):
16     def __init__(self, interval: timedelta, align_to: Optional[datetime] = None, init_item: Optional[int] = None):
17         super().__init__(interval, align_to)
18         logger.info(f"received starting id: {init_item}")
19         self.max_id = init_item
20
21     def next_item(self):
22         """
23         Get all the items from hacker news API between
24         the last max id and the current max id
25         """
26
27         if not self.max_id:

```

PROBLEMS 118 OUTPUT DEBUG CONSOLE TERMINAL

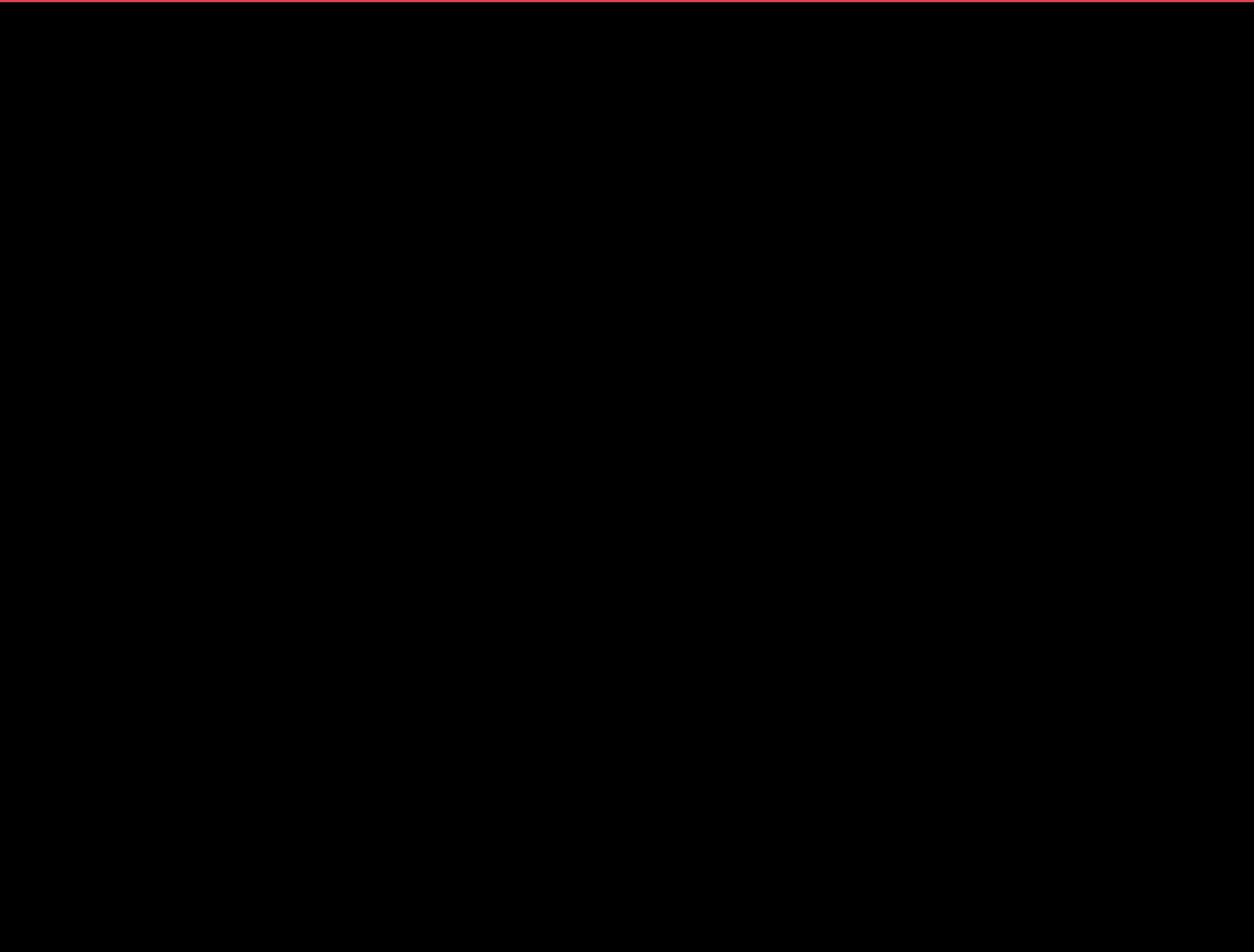
erefronts in two main areas, Russian-Ukraine and Israel-Palestine, seems to echo the situation before WWI and WWII. Not to mention coming at a time of intensifying geopolitical "chicken" and increased international distrust. Add that to a world suffering under the stress of economic crises, gas shortages and price hikes, and just out of a global pandemic—it seems as if we could be on a precipice. While there are some surface-level similarities between now and the pre-histories of both world wars of last century, there are also substantial differences. History doesn't repeat itself exactly, though it can sometimes rhyme. Differences include the presence of global institutions, nuclear deterrence, interconnected economies, the influence of 24/7 media, and a less ideological drive behind conflicts. Similarities feature regional conflicts that could draw in larger powers, economic crises and resource scarcity, a general atmosphere of global unrest, and rising nationalism. Alarmism is easy, and apocalyptic doomsayers are commonplace, so it's wise to be cautious and often sensible to dismiss such talk as baseless catastrophizing. However, it is possible that this time is different. What say you?', 'time': 1696764219, 'title': 'Ask HN: How close are we to a World War situation?', 'type': 'story'}}
<dataflow.Story object at 0x7f9f38162f80>
(('37805009', 1)
(('37809516', 4)
(('37809582', 2)
INFO:dataflow:current id: 37813690, new id: 37813690





Outline:

1. What and why real-time ML?
2. Building a pipeline to analyze streaming data
3. **Building a real-time feature pipeline**





Hacker News
Input

Add
Metadata

Add
Parent ID

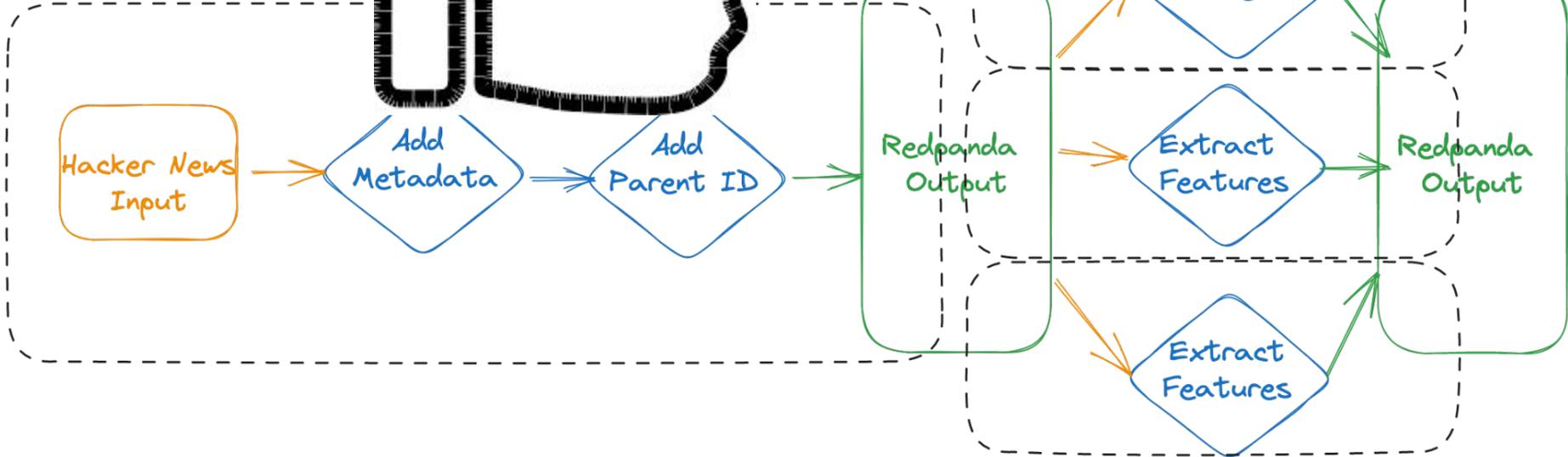
Redpanda
Output

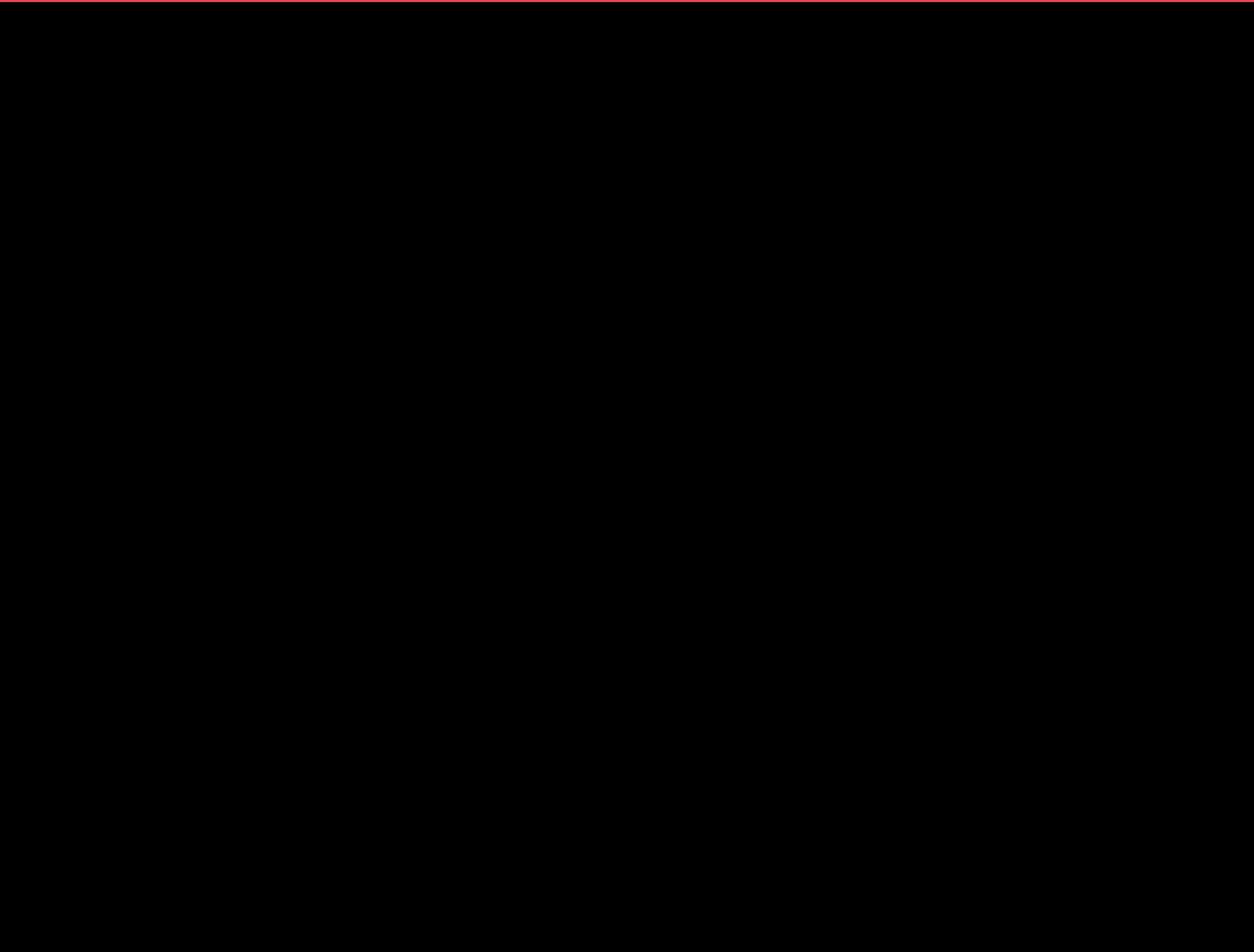
Count
Comments

Extract
Features

Redpanda
Output

Extract
Features







What is the next step to integrate this with feature stores for real-time ML?

With hopsworks you can register the features in the kafka topic and the ingestion into the hsfs will happen automatically. You can check out the tutorial in the logical clocks repo.

<https://github.com/logicalclocks/hopsworks-tutorials/tree/master/integrations/bytewax>



You can find the code used in the videos -

<https://github.com/awmatheson/guide-periodic-hackernews>

bytewax.io
github.com/bytewax/bytewax