# Handling Late Data and Watermarking
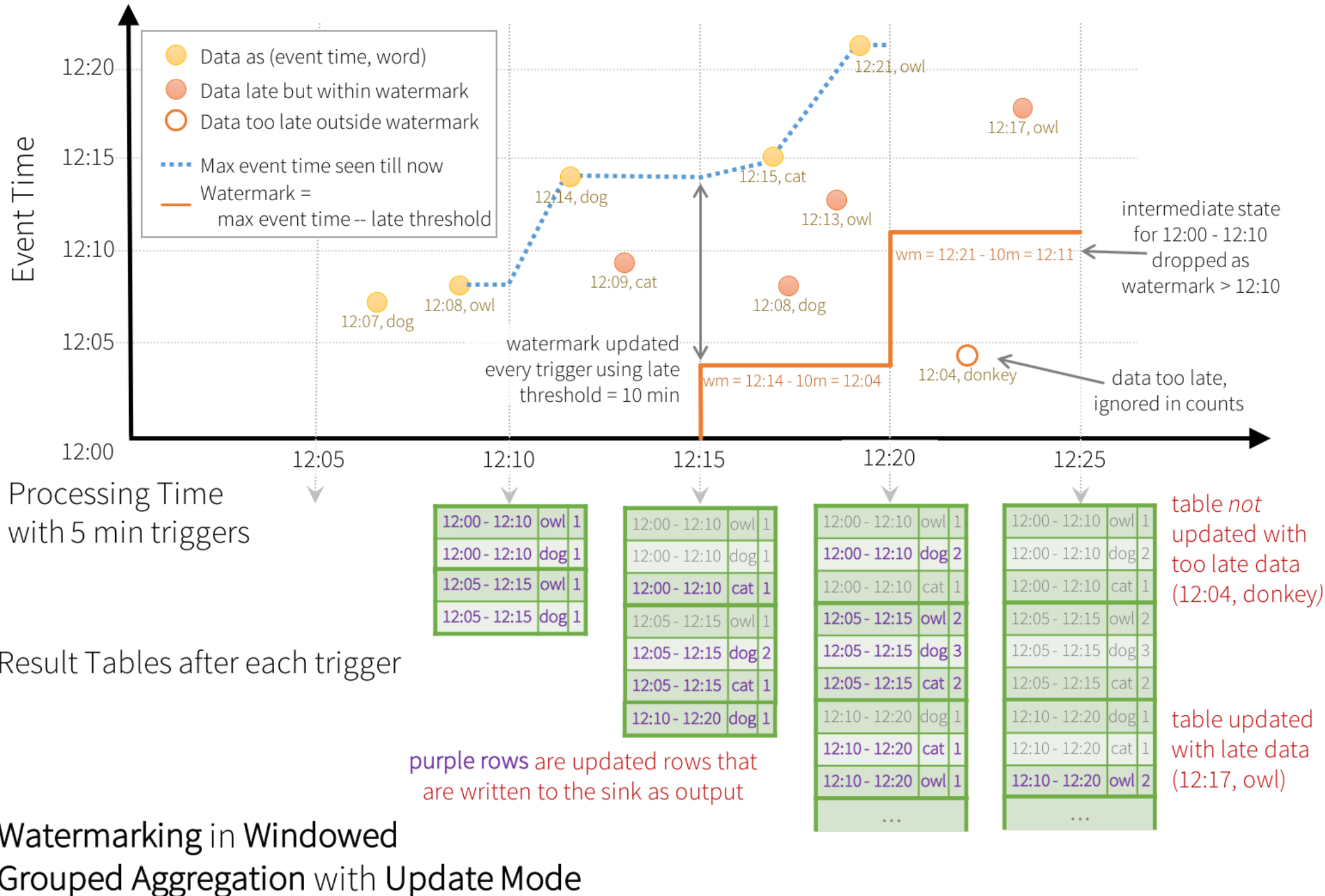
late data that was generated
at 12:04 but arrived at 12:11

Input Stream

| 12:02 | cat dog |
| 12:03 | dog dog |

| 12:07 | owl cat |

| 12:04 | dog |
| 12:13 | owl |

Time

12:00        12:05              12:10              12:15

Result Tables
after 5 minute triggers

| 12:00 - 12:10 | cat | 1 |
| 12:00 - 12:10 | dog | 3 |

| 12:00 - 12:10 | cat | 2 |
| 12:00 - 12:10 | dog | 3 |
| 12:00 - 12:10 | owl | 1 |
| 12:05 - 12:15 | cat | 1 |
| 12:05 - 12:15 | owl | 1 |

| 12:00 - 12:10 | cat | 2 |
| 12:00 - 12:10 | dog | 4 |
| 12:00 - 12:10 | owl | 1 |
| 12:05 - 12:15 | cat | 1 |
| 12:05 - 12:15 | owl | 2 |
| 12:10 - 12:20 | owl | 1 |

counts incremented only for
window 12:00 - 12:10

Late data handling in
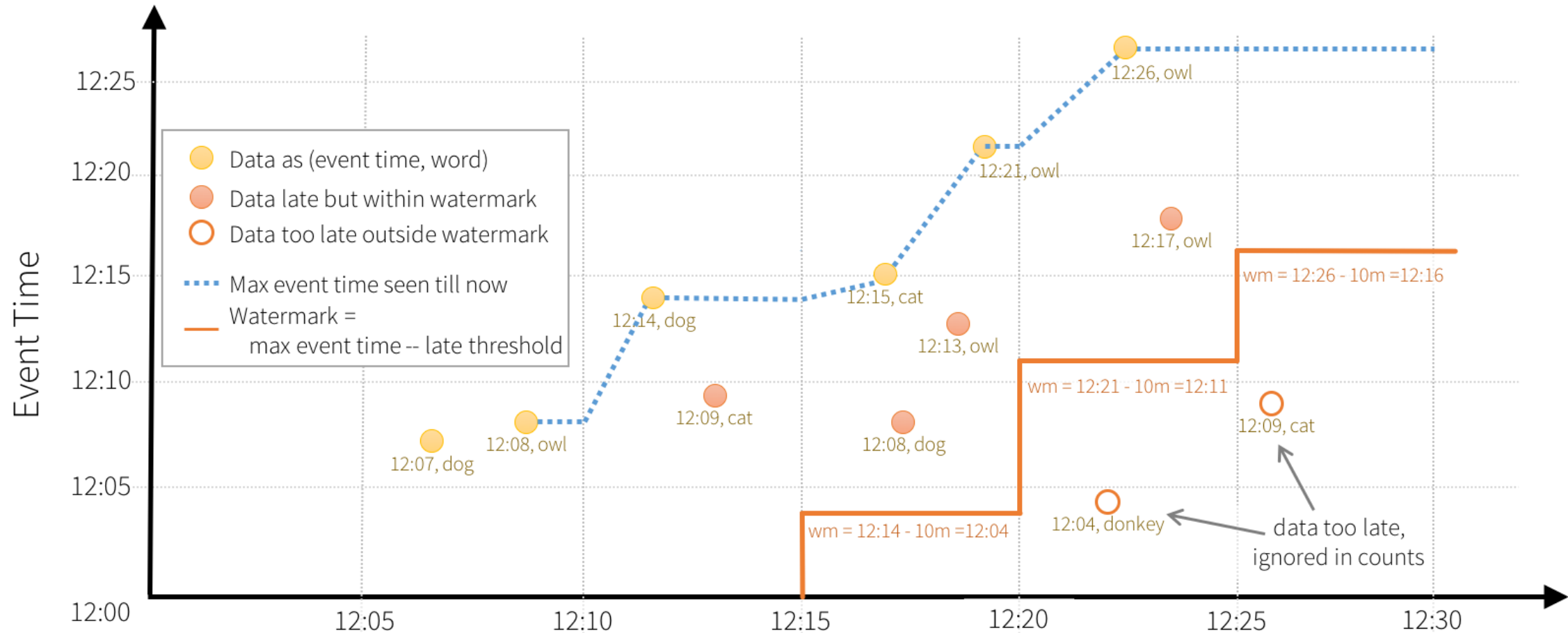Windowed Grouped Aggregation

# Watermarking example code

```python
words = ...   # streaming DataFrame of schema { timestamp:
Timestamp, word: String }


# Group the data by window and word and compute the count of
each group
windowedCounts = words \
    .withWatermark("timestamp", "10 minutes") \
    .groupBy(
        window(words.timestamp, "10 minutes", "5 minutes"),
        words.word) \
    .count()
```

Watermarking in Windowed Grouped Aggregation with Update Mode

Watermarking in Windowed Grouped Aggregation with Append Mode

Event Time

Processing Time with 5 min triggers

- Data as (event time, word)
- Data late but within watermark
- Data too late outside watermark
- Max event time seen till now
- Watermark = max event time -- late threshold

12:07, dog
12:08, owl
12:14, dog
12:09, cat
12:15, cat
12:13, owl
12:08, dog
12:21, owl
12:26, owl
12:17, owl
12:04, donkey
12:09, cat

wm = 12:14 - 10m = 12:04
wm = 12:21 - 10m = 12:11
wm = 12:26 - 10m = 12:16

data too late, ignored in counts

partial counts for window 12:00 - 12:10 maintained as internal state while waiting for late data, so not yet added to result table

final counts for 12:00 - 12:10 added to table when watermark > 12:10, late data counted, and intermediate state for window dropped

| 12:00 - 12:10 | owl | 1 |
| 12:00 - 12:10 | cat | 1 |
| 12:00 - 12:10 | dog | 2 |

| 12:00 - 12:10 | owl | 1 |
| 12:00 - 12:10 | cat | 1 |
| 12:00 - 12:10 | dog | 2 |
| 12:05 - 12:15 | owl | 2 |
| 12:05 - 12:15 | cat | 2 |
| 12:05 - 12:15 | dog | 3 |

Result Tables after each trigger

# Conditions for Watermarking

- Output mode must be set to _Append_ or _Update_.
  - If it is set to _Complete_ mode, it will requires all aggregate data to be preserved and will be incapable of using watermarking.

- Aggregation needs either an event-time column, or a window on the event-time column.
  - withWatermark must be called on the same column as the timestamp column used in the aggregate, `withWatermark` must be called before the aggregation for the watermark details to be used.