# Note on Zenoh

Hsin-Jui, Chen

December 12, 2024

## Contents

# 1 Tracing Zenoh

There are multiple ways to efficiently trace the Zenoh program using `uftrace`, one may follow these steps:

```
1 /* modify the toolchain */
2 $ cd $zenoh-repo/
3 $ vim rust-toolchain.toml                            // modify to "nightly"
4
5 /* install uftrace and compile zenoh */
6 $ sudo apt install uftrace
7 /* build zenoh using the following command lines */
8 $ export RUSTFLAGS=\  // Leading spaces are not allowed in the next line
9 "-Z instrument-mcount -C passes=ee-instrument<post-inline>"
10 $ cargo build --all-targets
11 $ uftrace record ./target/examples/z_pub
12 $ uftrace {replay,graph,tui}
```

To successfully compile the program and enable the program to traced in instruction-level. One may need to find out the right environment variables given to the compiler. I solve my problem with the comment in this reference, which provides some useful guidance for identifying the correct compile flag.

Another way to trace Zenoh programs is using the macros `tracing::info!()` and `tracing::debug!()` provided in crate `tracing`. These macros provide the same functionality as the `println!()` macro, however, it can separate the output into multiple levels of information. The following command is an example of running Zenoh program with `tracing` crate utilities.

```
1 $ RUST_LOG=debug cargo run --example example_name
2 /* You may also build and run */
3 $ cargo build --example
4 $ RUST_LOG=debug ./target/release/examples/example_name
```

To whom interested in how it works, you may check out the comment of the function `zenoh::try_init_log_from_env()` and its implementation. [1]

# 2 Case study: The `z_pub` example

```
0 # main() in examples/examples/z_pub.rs
30 let session = zenoh::open(config).await.unwrap();
```

The function `zenoh::open(config)` [2] returns an `OpenBuilder<TryIntoConfig>` instance, a structure that implements the `IntoFuture` trait. Therefore, the `.await` operator asynchronously waits for the `wait()` [3] function to complete.

---

[1] Line 31 in the file commons/zenoh-util/src/log.rs.
[2] Line 2900 in zenoh/src/api/session.rs
[3] Line 2954 in zenoh/src/api/session.rs

Then `wait()` function attempts to extract the config and passes it to the function `Session::new(config).wait()`[4]. This method first collects the topics it suscribes to and publishes. It then initializes a new runtime by

```
  0 # Session::new() in zenoh/src/api/session.rs
964 let mut runtime = RuntimeBuilder::new(config);
969 let mut runtime = runtime.build().await?;
```

Let's dive into the details of `runtime` and `RuntimeBuilder::build()`[5]. This function first sets up a router (responsible for managing the routing table), a transport manager (handling network protocols), and a notifier (notifying when subscribed topics are received). It then initializes the runtime state with the router, transport manager, and notifier. For reference, a brief snippet of the code is shown below:

```
  0 # RuntimeBuilder::build() in zenoh/src/net/runtime/mod.rs
  1 let router = Arc::new(Router::new(zid, whatami, hlc.clone(), &config)?);
143 let tm_builder = TransportManager::builder()
149     .from_config(&config)
150     .await?
151     .whatami(whatami)
152     .zid(zid);
159 let transport_manager = tm_builder.build(handler.clone())?;
169 let config = Notifier::new(config);
170 let runtime = Runtime {
171     state: Arc::new(RuntimeState {
185         ...                                              // ommited
186     }),
187 };
```

After that, it makes the process listen for incoming messages by starting an asynchronous notifier that listens for `"connect/endpoints"` with the runtime and updates peers as needed. However, I do not fully understand this part of the code, so I am only sketching the main ideas here. Once the runtime is built, the session is then initialized with the runtime.

```
  0 # Session::new() in zenoh/src/api/session.rs
971 let mut session = Self::init(
972     runtime.clone(),
973     aggregated_subscribers,
974     aggregated_publishers,
975 )
976 .await;
```

It appears that each session must be associated with a runtime; however, multiple sessions can share the same runtime. Within `Session::new()`, the `Session::init()`

---

[4] Line 955 in zenoh/src/api/session.rs
[5] Line 126 in zenoh/src/net/runtime/mod.rs

function is invoked, which subsequently calls `router.new_primitives()`[6]. This function adds a new `FaceState` entry to the routing table (`tables.faces`) and generates a unique `fid`. For reference, the relevant code snippet is shown below:

```
0    # Router::new_primitive() in zenoh/src/net/routing/router.rs
78   let fid = tables.face_counter;
79   tables.face_counter += 1;
80   let newface = tables
81       .faces
82       .entry(fid)
83       .or_insert_with(|| {
84           FaceState::new(
85               fid,
86               zid,
87               WhatAmI::Client,
88               #[cfg(feature = "stats")]
89               None,
90               primitives.clone(),
91               None,
92               None,
93               ctrl_lock.new_face(),
94           )
95       })
96       .clone();
97   tracing::debug!("New {}", newface);
98   for (key, val) in tables.faces.iter() {
99       tracing::debug!("key: {key} val: {val}");
100  }
```

From this, we can observe that the face ID (`fid`) is determined by a counter. Therefore, the observation mentioned in the meetings on 11/1 and 11/8, that `fid`s are inconsistent across Zenoh processes, is now clear and understandable.

After the `faces` is constructed, it is cloned using `Arc::clone()` and inserted into a `Face` structure, which is then returned. Subsequently, `admin::init(&session)` is called within `Self::init()`. However, the functionality of this method remains unclear to me. Once this step is complete, the session is successfully initialized and returned from `Self::init()` (`Session::init()`).

Next, the `Runtime::start()` function is invoked [7]. This function determines the appropriate method to call based on the value of runtime.WhatAmI and starts the runtime accordingly. In the `z_pub` example, the value of `WhatAmI` is set to `Peer`, resulting in the invocation of the `Runtime::start_peer()` method[8]. In the `Runtime` method `Runtime::start_peer()`, the following operations are performed:

---

[6] Implemented at line 70 in zenoh/src/net/routing/router.rs
[7] Line 119 in zenoh/src/net/runtime/orchestrator.rs
[8] Line 171 in zenoh/src/net/runtime/orchestrator.rs

**1.** Bind listeners (`self.bind_listeners(&listeners).await?;`): This step appears to create multiple listener threads to handle connections on different threads.

**2.** Connect to peers (`self.connect_peers(&peers, false).await?;`): Establishes connections to other peers.

**3.** Enable scouting (if configured): If scouting is enabled in the configuration, the runtime will start the scouting process.

When binding listeners, the critical functions in the call tree are as follows:

```
1 Runtime::bind_listeners(&listeners)
2   Runtime::bind_listeners_impl(listeners)
3     Runtime::add_listener(&self, listener: EndPoint)
4       Runtime::manager().add_listener(endpoint)
5         TransportManager::add_listener_unicast(endpoint)
6             TransportManager::new_link_manager_unicast(...)
7           LinkManagerUnicastTcp::new_listener(endpoint)
8             LinkManagerUnicastTcp::new_link_inner(...)
```

Then the TCP stream socket is returned.

# 3 On `io/zenoh-transport`

## 3.1 Basic Ideas of Transmission Pipeline

After some seniors advised me to switch my focus to studying the zenoh trasport, I start reading the codes in `io/zenoh-transport`, and I realized that this part might be more useful for understanding how Zenoh handles message transport.

I began by examining the code in `io/zenoh-transport/src/common/pipeline.rs`, which implements the pipeline transmission mechanism in Zenoh. The following note is written with respect to the earlier 1.0.0 beta version, it might have some differences now. This code establishes a MPSC (**Multiple Producer, Single Consumer**) pipeline. Consequently, the pipeline is divided into two stages: `StageIn` and `StageOut`.

For a data producer, it must insert data into a `WBatch` (short for *write batch*). The `WBatch` structure is defined in `io/zenoh-transport/src/common/batch.rs` as follows:

```
1 pub struct WBatch {
2     pub buffer: BBuf,              // The buffer to perform the batching on
3     pub codec: Zenoh080Batch,                    // The batch codec
4     /* It contains 1 byte to signal whether the batch is compressed */
5     pub config: BatchConfig,
6     #[cfg(feature = "stats")]          // Statistics related to this batch
7     pub stats: WBatchStats,
8 }
```

Essentially, the `WBatch` strucuture serves as an in-memory buffer for storing serialized data. To implement the MPSC pipeline, multiple `WBatch` structures are organized

into a `RingBuffer`. This design allows multiple producers to insert data into different batches concurrently. However, the implementation is more complex than it appears.

Next, we explain the key ideas behind the `push_network_message()` function in the `StageIn` implementation.

**1.** The function first attempts to obtain a `WBatch` by following these steps:

  **(a)** Acquire the lock associated with the `StageIn`.

  - The thread acquires a lock from `StageInMutex` to check if a `WBatch` is available.
  - If an available `WBatch` is found in the `StageInMutex`, the buffer is successfully obtained. The `WBatch` is then moved out of the `Option<WBatch>`, signaling to other threads that this `WBatch` needs to be refilled. [9]

  **(b)** Refill Check.

  - If no `WBatch` is available, the function attempts to refill by checking for remaining `WBatch` structures in the ring buffer. If successful, the thread obtains a buffer from the ring buffer.

  **(c)** Conditional Wait.

  - If the ring buffer also lacks an available `WBatch`, the thread must wait for a consumer to process a batch and notify the thread. This involves a conditional wait mechanism. Additionally, certain scenarios may require further considerations. For instance, if some messages are allowed to be dropped after a deadline, the thread may wait only until the deadline expires.

The following piece of code implements this idea:

```
1  macro_rules! zgetbatch_rets {
2    ($fragment:expr, $restore_sn:expr) => {
3      loop {
4        match c_guard.take() {
5          Some(batch) => break batch,
6          None => match self.s_ref.pull() {
7            Some(mut batch) => {
8              batch.clear();
9              self.s_out.atomic_backoff.first_write.store(
10                 LOCAL_EPOCH.elapsed().as_micros() as MicroSeconds,
11                 Ordering::Relaxed);
12              break batch;
13            }
14            None => {
15              drop(c_guard);
16              match deadline_before_drop {
```

---

[9] Initially, I thought that moving the `WBatch` out of the `Option<T>` was to prevent other threads from writing to the same `WBatch`. Thanks to those seniors for pointing out my misunderstanding.

```
17                        Some(deadline) if !$fragment => {
18                            if !self.s_ref.wait_deadline(deadline) {
19                                $restore_sn;
20                                return false
21                            }
22                        }
23                        _ => {
24                            if !self.s_ref.wait() {
25                                $restore_sn;
26                                return false;
27                            }
28                        }
29                    }
30                    c_guard = self.mutex.current();
31                }
32            },
33        }
34    }
35    };
36 }
```

**2.** After obtaining the buffer, the function attempts to write the message into it. However, the acquired `WBatch` may not be empty, and the message might not fit into the remaining space. In such cases, the function releases the original `WBatch` and tries to obtain an empty `WBatch` using the same procedure described above.

**3.** If an empty `WBatch` still cannot accommodate the entire message, the function will iteratively obtain additional `WBatch` structures and fragment the message until it is fully processed.

Note that there is a distinction between `NetworkMessage` and `TransportMessage`: the former represents the data being transported, while the latter is used for protocol and flow control. As a result, a network message may be larger and may not fit into a single `WBatch`. However, in `push_transport_message()`, fragmentation handling is unnecessary. To implement the concepts described above, the following structures are designed to manage the `StageIn` and `StageOut` pipelines.

```
1 /* StageIn */
2 struct StageInRefill {
3     n_ref_r: Waiter,
4     s_ref_r: RingBufferReader<WBatch, RBLEN>,
5 }
6 struct StageInOut {
7     n_out_w: Notifier,
8     s_out_w: RingBufferWriter<WBatch, RBLEN>,
9     atomic_backoff: Arc<AtomicBackoff>,
10 }
```

```
11  struct StageInMutex {
12      current: Arc<Mutex<Option<WBatch>>>,
13      priority: TransportPriorityTx,
14  }
15  /* StageOut */
16  struct StageOutIn {
17      s_out_r: RingBufferReader<WBatch, RBLEN>,
18      current: Arc<Mutex<Option<WBatch>>>,
19      backoff: Backoff,
20  }
21  struct StageOutRefill {
22      n_ref_w: Notifier,
23      s_ref_w: RingBufferWriter<WBatch, RBLEN>,
24  }
```

We can observe that the `StageInRefill` structure holds the information, methods, and implementations necessary for waiting on a `WBatch` to be consumed. On the other hand, the `StateInOut` structure stores the information required to notify the consumer threads.

## 3.2   Backoff Mechanism in Transmission Pipeline

As I mentioned in previous meeting, there are some kind of backoff mechanisms in the transmission pipeline. I believed that it is for conjestion control, however, as I studied more and more, I find that this is totally not true.

The idea is that we do not want a `WBatch` be consumed if it only have a few bytes written. Hence, the backoff is adopted to avoid constantly writing samll batches. That is, if only a few bytes are written to a `WBatch`, then there is no need to consume it instantly. In this case, waiting for some bytes begin produced is propably making the transmission more efficient and less communication overhead (like header).

I will elaborate more details here. Recall that both consumer and producer maintain their availble `WBatch` queues. For consumer, the queue (ring buffer) contains the `WBatchs` ready to send out. On the other hand, the queue for producer contains the `WBatch` that is empty (but it might not be clean). The advantage for using two queues is that it limits the maximum waiting `WBatch` and can be viewed as a resource controller. For consumer, if there is a `WBatch` in its queue, then it just pull it and send out on the link. Clearly, if a `WBatch` is in the consumer queue, it indicates that it either full or the message need to be send as soon as possible. In this scenario, there is no need to wait.

However, when the consumer queue is empty, the consumer might want to consume the currently processing `WBatch` in `StageIn`. To prevent often sending out `WBatch` with only few bytes, we make the consumer wait until one of the conditions are met:

1. A `WBatch` is put into the consumer queue. This will be notified by producers.

2. The `WBatch` in `StageInMutex` contains enough bytes (over a threshold). This will

be notified by producers.

3. The backoff timeout.

In the first case, it means that a WBatch becomes ready when waiting. In the second case, it means that a WBatch is nearly full and can be send out. The third case indicates that no new messages are written to the current WBatch and there is no need to wait longer. We now look at the code and give some explanations.

```
1 struct AtomicBackoff {
2     active: CachePadded<AtomicBool>,
3     bytes: CachePadded<AtomicBatchSize>,
4     first_write: CachePadded<AtomicMicroSeconds>,
5 }
```

This structure is shared by producer and consumer. Three fields in the structure has the following meaning:

1. active indicates whether the consumer is actively backoff.

2. bytes stores the number of bytes written in the current WBatch.

3. first_write stores the time that the first message written to current WBatch. This is used to prevent a message waits for too long.

The following code show how the ideas are implemented in producer side.

```
1 impl StageInOut {
2     #[inline]
3     fn notify(&self, bytes: BatchSize) {
4         self.atomic_backoff.bytes.store(bytes, Ordering::Relaxed);
5         if !self.atomic_backoff.active.load(Ordering::Relaxed) {
6             let _ = self.n_out_w.notify();
7         }
8     }
9     #[inline]
10     fn move_batch(&mut self, batch: WBatch) {
11         let _ = self.s_out_w.push(batch);
12         self.atomic_backoff.bytes.store(0, Ordering::Relaxed);
13         let _ = self.n_out_w.notify();
14     }
15 }
```

In the function notify(), it will first store the number of bytes written, and if the consumer is backoffing (the active field is true), then it will notify the consumer and the consumer will be waken up.

For the function move_batch(), it will push the WBatch to the consumer queue, and since the current WBatch is moving out, we need to set the number of bytes to 0. Afterwards, it will notify the consumer.

On the consumer side, two methods `try_pull()` and `try_pull_deep()` of the structure `StageOutIn` implement the functionality of pulling message from the queue or the current `WBatch`. The following pseudo code demonstrate the workflow of pulling `WBatch` from consumer.

---

**Algorithm 1** The implementations of `try_pull()` and `try_pull_deep()` functions

---

1: **function** TRYPULL( )
2:     **if** the consumer queue is not empty **then**
3:         set the backoff state to inactive         *// it might be backoffing previously.*
4:         **return** The `WBatch` in the queue
5:     **else**
6:         TRYPULLDEEP( )         *// Try to pull current WBatch.*
7:
8: **function** TRYPULLDEEP( )
9:     *pull* ← (backoff is active)     *// Check whether it was already in backoff since last pull*
10:     *backoffTime* ← 0
11:     **if** not *pull* **then**         *// The consumer is already backoffed in previous pulls*
12:         *pull*←(new_bytes == old_Bytes)     *// Check whether some new bytes are added*
13:     **if** not *pull* **then**         *// There are some bytes added by the producer*
14:         *diff* ← elapsed time since the first write to the `WBatch`
15:         **if** *diff* ≥ *threshold* **then**
16:             *pull* ← True
17:         **else**
18:             *backoffTime* ← *threshold - diff*
19:     **if** *pull* **then**         *// Ready to pull the current WBatch.*
20:         acquire lock
21:         set backoff state to inactive
22:         **if** consumer queue is not empty **then**     *// In case the queue is not empty*
23:             **return** `WBatch` from the queue
24:         **else if** current `WBatch` does not exist **then**     *// StageIn does not refill yet*
25:             **return** None
26:         **else**         *// Pull the current WBatch*
27:             **return** current `WBatch`
28:     **else**         *// The requirements of pulling current WBatch does not meet*
29:         set backoff state to active
30:         **return** *backoffTime*     *// Make consumer wait for some time*

---

This makes me realize some of my previous ideas on transmission pipeline is actually incorrect, and may have some subtle (but maybe neglectable) impact on my later `qstats` feature implementations for per-pipeline profiling.

For example, a statistic we may be interested in is the average queueing delay of messages or `WBatches`. Suppose we are profiling `WBatches`' queueing dealy, intuitively, we may record the time when a `WBatch` be put into the queue, then find the elapsed time when the same `WBatch` is being pulled. However, consider a scenario that a message has been written to a new `WBatch`, and that queue has not been consumed at all,

then this `WBatch` contains some data but never start counting the queueing delay, which somehow feels not right.

# 4 Profiling Zenoh

## 4.1 Official Implementation of the `stats` Feature

In this section, I study the file `io/zenoh-transport/src/common/stats.rs`, and I am going to briefly expalin how Zenoh team implement their statistics analysis tools.

Before we get started, I will first show how we can utilize these existing implementations. First, to enable this, you may need to enable the `stats` features when compiling, you may do it in many ways, for instance, the following compiling commands all work.

```
1 $ cargo run --example <example_name> --features stats
2 /* it is also possible separately build and run */
3 $ cargo build --example <example_name> --features stats
4 $ ./target/debug/<example_name>
```

To get famaliar with the `TransportStats` structure, I first implement printing out the statistics on dropping, that is, as soon as the `TransportUnicast` drop, then it will print out the statistics. Precisely speaking, I implement the `Drop` trait for the structure `TransportUnicastUniversal` as the following code listing. [10]

```
1 impl Drop for TransportUnicastUniversal {
2     fn drop(&mut self) {
3         tracing::info!(
4             "{}",
5             self.stats.clone().report().openmetrics_text()
6         );
7     }
8 }
```

Since the implementation is written in a complicated macro, I decide to expand the macro to have better understand on this complicated macro. I use the following commands to expand the macro.

```
1 $ cargo install cargo-expand
2 $ cd io/zenoh-transport/src
3 $ cargo expand common::stats --features stats > common/stats_expand.rs
```

We can see that it actually implement some increase counter methods and getter methods, this idea gives me some insights to build the equivalent of `qstats`.

## 4.2 Documentation for the `qstats` Feature

This subsection I am trying to build the `qstats` feature by imitating the `stats` feature after macro expansions. The code is available at this repository. I give definitions

---

[10] I implement this trait in the file io/zenoh-transport/src/unicast/universal/transport.rs.

and implement some methods of the structures `QueueStats` and `QueuesStatsReport` in the file `io/zenoh-transport/src/common/qstats.rs`.

```rust
pub struct QueueStatsReport {
    priority: usize,
    pub avg_qsize: f64,
    pub droprate: f64,
    pub avg_qdelay: f64,
    pub mcnt: usize,
}


pub struct QueueStats {
    queue_counter: AtomicUsize,   // current number of queueing messages
    priority: usize,                         // the priority of this queue
    qsize: Arc<Mutex<Vec<usize>>>,       // recorded numbers of messages
    dropped: AtomicUsize,                // the number of dropped messages
    tried: AtomicUsize,                   // the number of tried messages
    queueing_delay: Arc<Mutex<Vec<usize>>>, //  recorded queueing delays
}
```

These are the definitions of the structures, and we can see that `QueueStatsReport` will record the `QueueStats` at some point, while `QueueStats` structure continuously record the statistics. The main statistics we are concerned about are average queue size (`qsize`), number of dropped messages (`dropped`), drop rate due to queueing (`dropped/tried`), and average queueing delay (average of `queueing_delay`).

The structure `QueueStats` have the following methods:

```rust
impl QueueStats {
    pub fn new(priority: usize) -> Self {}
    pub fn record_qsize(&self) {}
    pub fn push_qdelay(&self, delay: usize) {}
    pub fn inc_qcnt(&self) {}
    pub fn dec_qcnt(&self) {}
    pub fn inc_dropped(&self) {}
    pub fn inc_tried(&self) {}
    pub fn report(&self) -> QueueStatsReport {}
}
```

These methods have the following purposes:
1. `new()`: create a new `QueueStats` for priority level `prioirty`.
2. `record_qsize()`: when producer trying to push message, it will first need to record the number of messages in the queue. Hence the method need to first load the counter then push to `self.qsize`.
3. `push_qdelay()`: when a consumer pull a message, it will calculate the elapsed time since the message written to a `WBatch`.
4. `inc_qcnt()`, `dec_qcnt()` when a message is pushed or pulled, we need to modify

the counter.

5. `inc_dropped()`, `inc_tried()` when a message is pushed, we need to change the counter of dropped and tried messages.

6. `report()` generate a `QueueStatsReport` structure that contains the current statistics. For the average queue size and average queueing delay, I choose to implement the exponential moving average formula, with $\alpha = 0.1$.

**Modifications on `batch.rs`.** There are modifications needed to be made to `batch.rs` and `pipeline.rs`. For `batch.rs`, we add the following to `encode()` [11] methods. We use the conditional compilation flag `#[cfg(feature = "qstats")]` to ensure the modifications are effective only when compiling with qstats feature.

```
 1 struct WBatch {
 2     /* some fields are omitted */
 3     #[cfg(feature = "qstats")]
 4     pub time: Vec<Instant>,
 5 }
 6 impl Encode<...> for &mut WBatch {
 7     fn encode(...) -> Self::Output {
 8         /* some lines are omitted */
 9         let res = self.codec.write(&mut writer, x);
10         #[cfg(feature = "qstats")]
11         if res.is_ok() {
12             self.time.push(Instant::now());
13         }
14         res
15     }
16 }
```

When a message or a fragment buffer is successfully written to a `WBatch`, we add an instant (time) to `WBatch.time`. Then, when this `WBatch` is pulled, it will calculate queueing delay of these messages and fragments.

**Modifications on `pipeline.rs`.** In this paragraph, we only list important modifications. The consumer and the producer both will have a list of `QueueStats`.

```
 1 struct StageIn {
 2     /* some fields are omitted */
 3     #[cfg(feature = "qstats")]
 4     qstats: Arc<QueueStats>,
 5 }
 6 pub(crate) struct TransmissionPipelineProducer {
 7     /* some fields are omitted */
 8     #[cfg(feature = "qstats")]
 9     qstats_list: Arc<[Arc<QueueStats>]>,
10 }
```

---

[11] Note that there are multiple implementations of `encode()`.

```rust
11 pub(crate) struct TransmissionPipelineConsumer {
12     /* some fields are omitted */
13     #[cfg(feature = "qstats")]
14     qstats_list: Arc<[Arc<QueueStats>]>,
15 }
```

Although the `QueueStats` structure for each queue will be access by consumer and producer simultaneously, it is not needed to use `Mutex` because the inner data structure of `QueueStats` is protected by atomic operations and mutexes. That is, the interfaces provide by `QueueStats` are already thread-safe. With the modifications above, we can now look at the actual changes in the consumer and the producer.

```rust
1  impl TransmissionPipelineConsumer {
2    pub(crate) async fn pull(&mut self) -> Option<(WBatch, usize)> {
3      while self.active.load(Ordering::Relaxed) {
4        let mut backoff = MicroSeconds::MAX;      // Calculate the backoff
5        for (prio, queue) in self.stage_out.iter_mut().enumerate() {
6          match queue.try_pull() {
7            Pull::Some(batch) => {
8              #[cfg(feature = "qstats")]
9              {
10               let qstat = &self.qstats_list[prio];
11               batch.time.iter().for_each(|t| {
12                 qstat.dec_qcnt();
13                 qstat.push_qdelay(t.elapsed().as_micros() as usize);
14               });
15             }
16             return Some((batch, prio));
17           }
18           Pull::Backoff(deadline) => {
19             backoff = deadline;
20             break;
21           }
22           Pull::None => {}
23         }
24       }
25       /* the remaining part is omitted */
26     }
27   }
28 }
```

We see that in the conditional compile block, for each `Instant` on the `WBatch`, we will calculate the delay and decrease the counter by the number of messages or fragments. The next part the changes on producer side.

```rust
impl TransmissionPipelineProducer {
    pub(crate) fn push_{network,transport}_message(
        &self,
        mut msg: NetworkMessage
    ) -> bool {
        /* some lines are ommited */
        let result = queue.push_network_message(
            &mut msg, priority,
            &mut deadline
        );
        #[cfg(feature = "qstats")]
        {
            let qstats = &self.qstats_list[idx];
            qstats.record_qsize();
            qstats.inc_tried();
            if !result {
                qstats.inc_dropped();
            }
        }
        result
    }
}
```

In this part, we need to record the queue size when a message is pushed and increase the counter of tried messages. If the result is not successful, the counter of dropped messages should also increase. The counter of queueing messages is added in the `zretok!()` macro defined in `StageIn` methods `push_network_message()` and `push_transport_message()` and the fragment part of `push_network_message()`.

## 4.3  Documentation for the Link Scheduler

In this subsection, I am going to explained how I abstracted the original scheduler into `Scheduler` trait.

```rust
trait SchedulerTrait {
    fn reset(&mut self);
    fn schedule(&mut self) -> ScheduleResult;
}
enum ScheduleResult {
    Some((WBatch, usize)),
    Backoff(MicroSeconds),
    None,
}
```

The above piece of code defines the `Scheduler` trait and the `SchedulerResult` enumeration, which is the result returned by `scheduler()` function. The following structure is the structure that implements the original schedule algorithm:

```rust
1  struct PriorityScheduler {
2      stage_out: Arc<Mutex<Box<[StageOut]>>>,
3  }
4  impl SchedulerTrait for PriorityScheduler {
5      fn reset(&mut self) -> {}                            // just do nothing
6      fn schedule(&mut self) -> ScheduleResult {
7          let mut stage_out = zlock!(self.stage_out);
8          for (prio, queue) in stage_out.iter_mut().enumerate() {
9              match queue.try_pull() {
10                 Pull::Some(batch) => {
11                     return ScheduleResult::Some((batch, prio));
12                 }
13                 Pull::Backoff(deadline) => {
14                     return ScheduleResult::Backoff(deadline);
15                 }
16                 Pull::None => {}
17             }
18         }
19         ScheduleResult::None
20     }
21 }
```

Basically, I just copy the original code into this function. Here, I am going to give another example of scheduler. I implemented a scheduler that combined weighted round robin and priority scheduling. The idea is that we divide the transmission into two phases: contention and polling. In the contention phase, the scheduler will schedule based on priority, while in the polling phases, the scheduler will based on the allocated virtual time.

```rust
1  struct WRRScheduler {
2      stage_out: Arc<Mutex<Box<[StageOut]>>>,
3      c_window_start: Option<Instant>,
4      n_window_start: Option<Instant>,
5      last_update: Option<(Instant, MicroSeconds, usize)>,
6      virtual_time: BinaryHeap<Reverse<(MicroSeconds, usize)>>,
7      weight: Vec<usize>,
8      c_window: MicroSeconds,
9      n_window: MicroSeconds,
10     is_contention: bool,
11     length: usize,
12 }
```

Some fields store the setting of the scheduler, for instance, c_window and n_window represent the length of the contention window and the polling window. The weight represent the weight of all priorities. Unfortunately, the explanation of the scheduler described above is temporarily unavailable. I am still working on it and hope to resolve this soon.

After having implemented our first customized scheduler, we may also define the enumeration `SchedulerConfig` and the structure `SchedulerBuilder`. The structure `SchedulerBuilder` implements the `build()` method that will return an instance of `Box<dyn SchedulerTrait + Send + Sync>` generics. Precisely speaking, the `build()` method will return the `Scheduler` trait object that is specific by the configuration. [12]

```rust
#[derive(Default)]
#[allow(dead_code)]
enum SchedulerConfig {
    #[default]
    Priority,
    WeightedRoundRobin(WRRConfig),
    Custom(CustomConfig),
}
struct SchedulerBuilder {
    config: SchedulerConfig,
}
```

With those enumeration and structure defined, we can implement the `build()` and other related methods as follow.

```rust
impl SchedulerBuilder {
    fn build(
        self,
        stage_out: Arc<Mutex<Box<[StageOut]>>>,
    ) -> Box<dyn SchedulerTrait + Send + Sync> {
        match &self.config {
            SchedulerConfig::Priority => self.build_priority(stage_out),
            SchedulerConfig::WeightedRoundRobin(config) => {
                self.build_wrr(stage_out, config.clone())
            }
            SchedulerConfig::Custom(config) => {
                self.build_custom(stage_out, config.clone())
            }
        }
    }
    fn build_priority(
        &self,
        stage_out: Arc<Mutex<Box<[StageOut]>>>
    ) -> Box<PriorityScheduler> {
        Box::new(PriorityScheduler { stage_out })
    }
    fn build_wrr(
        &self,
        stage_out: Arc<Mutex<Box<[StageOut]>>>,
        config: WRRConfig,
```

---

[12] Send and Sync is added to make compiler not complaining about multi-threading.

```
26        ) -> Box<WRRScheduler> {
27            let length = stage_out.lock().unwrap().len();
28            Box::new(WRRScheduler {
29                stage_out,
30                c_window_start: None,
31                n_window_start: None,
32                last_update: None,
33                virtual_time: BinaryHeap::new(),
34                weight: config.weight,
35                c_window: config.c_window,
36                n_window: config.n_window,
37                is_contention: false,
38                length,
39            })
40        }
41        /* fn build_custom() is omitted */
42 }
```

First, the build() method will call the corresponding sub-build method based on the given configuration. Next, the sub-build method will the construct the scheduler structure then returned.