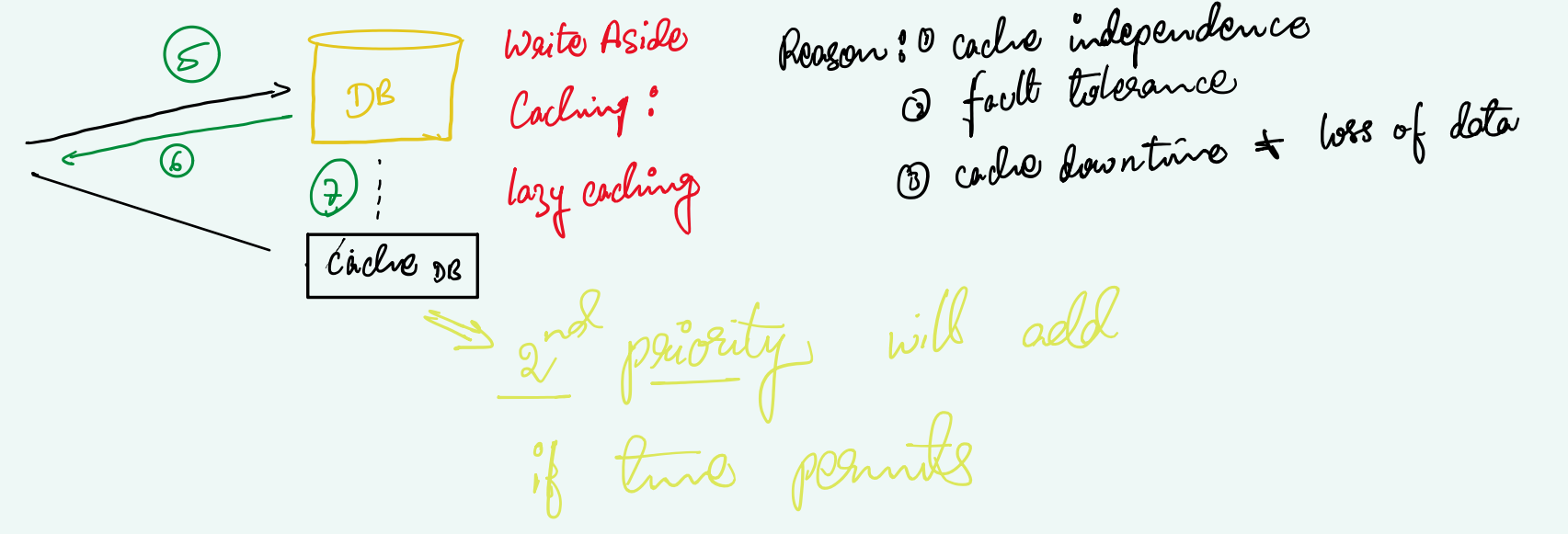


localhost: 8081 / endpoint.



Ingestion Service
middleware check
data structure validation
events not data prep
ETL processing (if req)
DB writes



Response is independent of successful DB storage

Response types
① 200 OK - successfully ingested

② error response in case of failure in ingestion

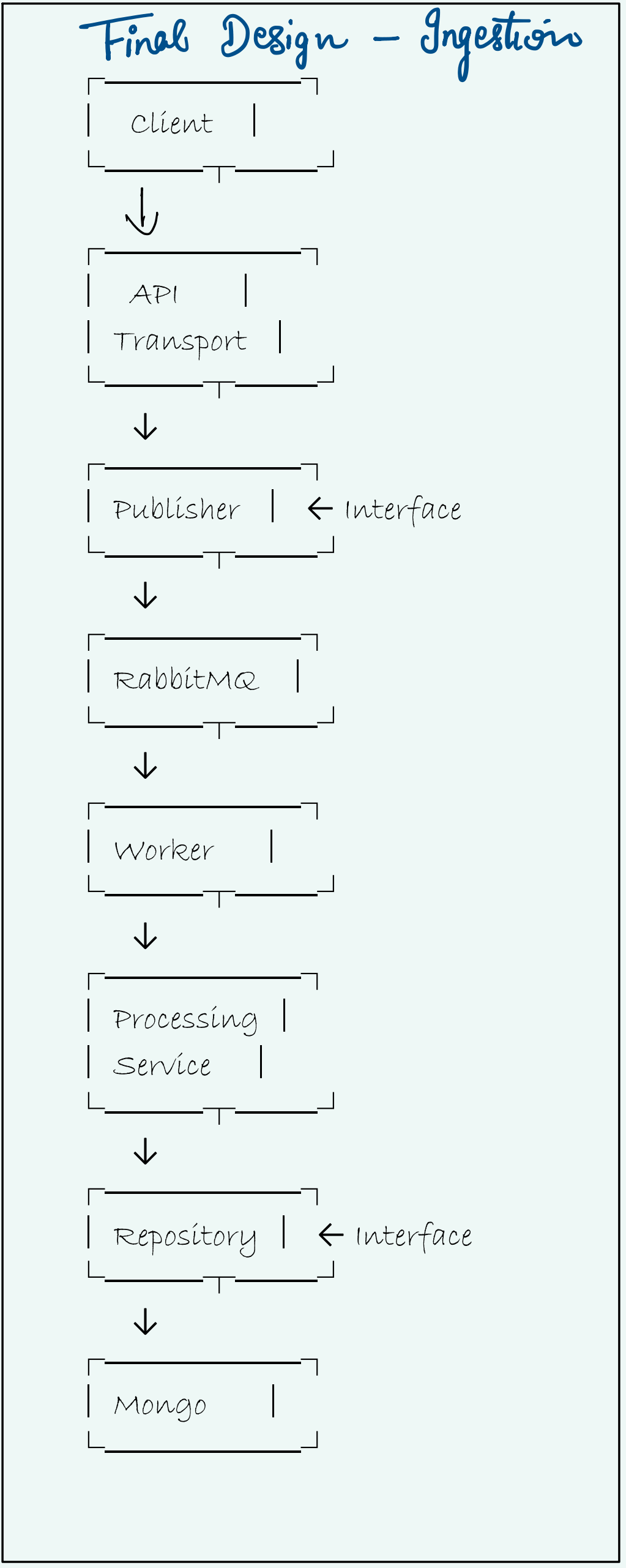
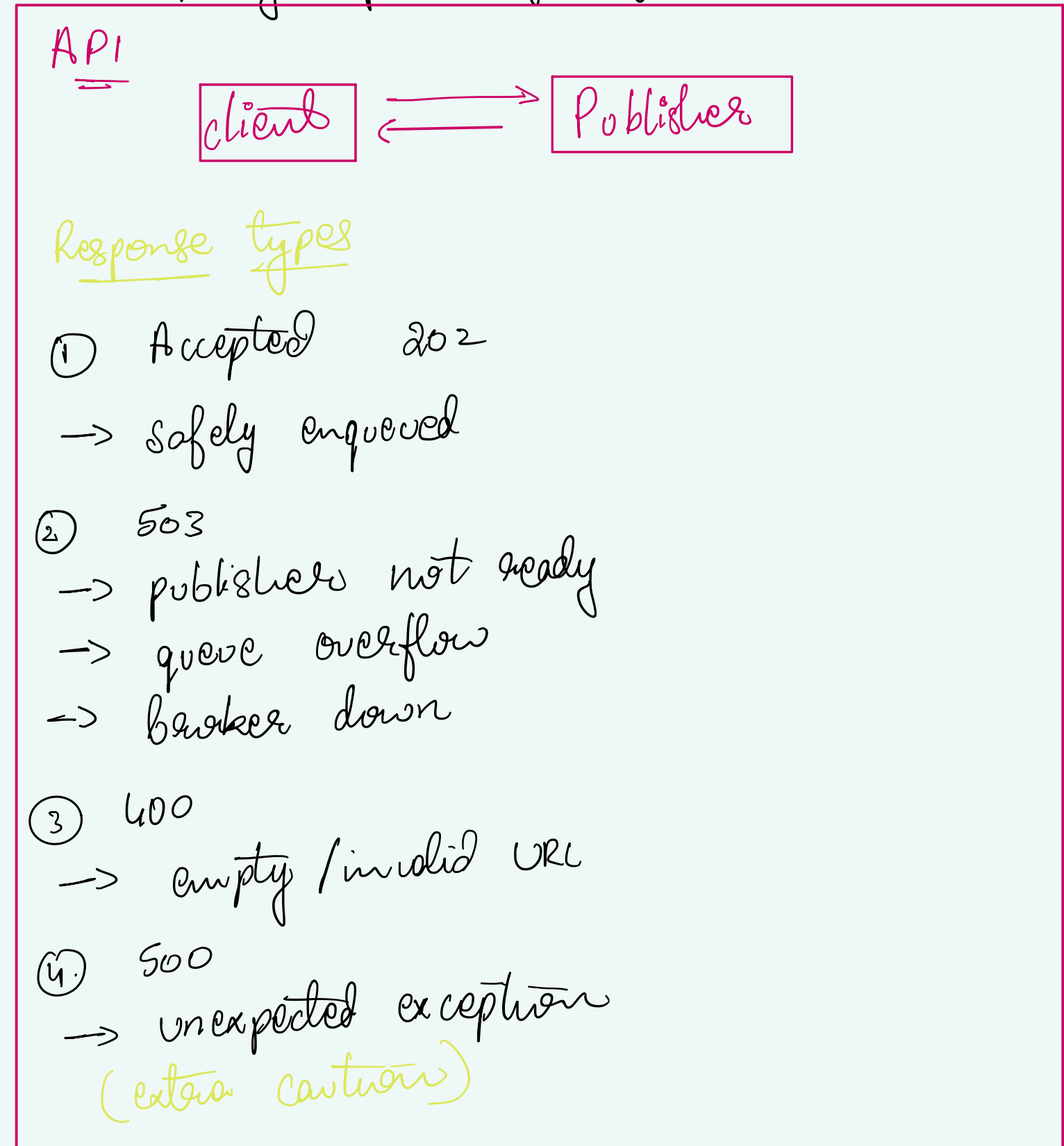
eg: capacity of queue is reached.
→ queue unavailability

Assumptions

- ① Data Consistency is of prime importance
- ② Data Availability is 2nd priority {older data might be fetched from cache.}
- ③ A cache invalidation API can be setup but not advisable post every DB update
- ④ Cache TTL can be set to 1 hour {ensuring everything before that is updated correctly.}

4 Sep Containers in Docker.

- ① API
 - * health endpoints
 - * post
 - * get
 - * exp backoff retries to connect with queue & DB
- ② worker
 - * exp backoff retries to connect with queue & DB
 - * manage states
 - * consume msg
 - * retry logic
 - * update in DB
- ③ DB
 - MongoDB
 - volume mount
- ④ RabbitMQ
 - * Durable
 - * persistent
 - * volume mounted (disk)
 - * prefetch = 1
 - * max length = 1000
 - * Default priority (for MVP no retry ④)
 - * reject publish if overflow



Design Schema (Booq)

```
{
  _id: ObjectId,
  url: string, // unique
  status: string,
  headers: object | null,
  cookies: object | null,
  page_source: string | null,
  error_message: string | null,
  retry_count: int,
  created_at: datetime,
  updated_at: datetime,
  last_attempt_at: datetime | null
}
```

cons: ① model modifications would be needed in future if any new data needs to be recorded
② not extensible
③ too many fields in model.

DB Schema design

Indexes:
① URL *** for faster query
② created_at ⇒ for querying in date range

Final Schema

```
{
  "_id": ObjectId,
  "url": "string",
  "status": "QUEUED | IN_PROGRESS | COMPLETED | FAILED_RETRYABLE | FAILED_PERMANENT",
  "metadata": {
    "headers": {...} | null,
    "cookies": {...} | null,
    "page_source": "string" | null
  },
  "processing": {
    "retry_count": int,
    "error_message": "string" | null,
    "last_attempt_at": datetime | null
  },
  "additional_details": {...} | null,
  "created_at": datetime,
  "updated_at": datetime
}
```

Processing Flow

On message receive:

- Upsert record (if not exists)
- Transition → IN_PROGRESS
- Fetch metadata
- If success:
 - Update status → COMPLETED
 - Store headers, cookies, page_source
 - ACK
- If failure:
 - Increment retry_count
 - If retry_count < 3:
 - Update status → FAILED_RETRYABLE
 - NACK (requeue)
 - Else:
 - Update status → FAILED_PERMANENT
 - ACK