



# Review of Pipeline Documentation and Code

## Naming Policy Violations

The repository contains several instances where file naming does not conform to the prescribed naming conventions <sup>1</sup>. According to the policy, pipeline documentation files should be named as `NN-<entity>-<provider>-<topic>.md` and pipeline code modules as `<stage>.py` under the `src/bioetl/pipelines/<provider>/<entity>/` path <sup>2</sup>. The following violations were identified:

- **Documentation Files:**

- **Target pipeline:** `docs/02-pipelines/chembl/target/01-target-normalizer.md` - Missing provider name in filename (should include "chembl") <sup>3</sup>. This violates **DOC-004** which requires the pattern `<entity>-<provider>-<topic>` <sup>1</sup>.
- **Assay pipeline:** `docs/02-pipelines/chembl/assay/02-assay-normalizer.md` and `03-assay-payload-parser.md` - Missing provider name in filenames <sup>4</sup>. They should be named `02-assay-chembl-normalizer.md` and `03-assay-chembl-payload-parser.md` to include the provider.
- **TestItem pipeline:** Multiple files in `docs/02-pipelines/chembl/testitem/` deviate from the naming scheme. For example, `02-chembl-testitem-thin-pipeline.md` has the provider and entity order swapped (should be `02-testitem-chembl-thin-pipeline.md`), and files like `03-testitem-pipeline-options.md`, `04-parent-enrichment-preparation.md`, etc., omit the provider/entity entirely <sup>5</sup>. All pipeline docs must include both entity and provider in the filename per the naming policy.

- **Common component docs:** Under `docs/02-pipelines/chembl/common/`, some files do not follow a consistent pattern. For instance, `14-chembl-client.md` has only two segments after the number (violating the `NN-entity-provider-topic` format), and files like `23-extraction-strategy-factory.md`, `24-service-extraction-strategy.md`, `25-dataclass-extraction-strategy.md` lack the provider segment <sup>6</sup>. Conversely, others like `18-chembl-base-pipeline.md` and `20-chembl-extraction-service.md` include the provider as the first segment <sup>7</sup>, which is the reverse order. These inconsistencies violate the deterministic naming rule **DOC-004** <sup>1</sup>. All documentation filenames should be corrected to match the `<entity>-<provider>-<topic>` format (or an approved variant if "common" is treated specially).

- **Code Files:**

- **Aggregated stage modules:** In the pipeline code, modules grouping multiple stages (e.g. `src/bioetl/pipelines/chembl/activity/stages.py`) break the intended one-stage-per-module pattern. The naming policy requires individual stage files such as `extract.py`, `transform.py`, `validate.py`, etc., under each pipeline directory <sup>2</sup> <sup>8</sup>. A file like `stages.py` or a generically named `etl.py` is considered a violation (monolithic pipeline module, flagged as **PIPE-004**) <sup>9</sup> <sup>10</sup>. For example, the ChEMBL Activity pipeline currently uses `stages.py` containing `ActivityExtractor/Transformer/Writer` classes, whereas it **should** be split into `extract.py`, `transform.py`, `write.py` for clarity and compliance.

- **Incorrect stage naming:** Any stage file names that use nouns or plurals instead of the verb form should be renamed. The policy dictates using action names (e.g. `validate.py` instead of `validator.py`)<sup>8</sup>. In this codebase, the ChEMBL Activity pipeline previously had `validator.py` (as noted in the style guide) which **should** be `validate.py`<sup>8</sup>. Similarly, `normalizers.py` (e.g. for assay and target pipelines<sup>11</sup>) would be more consistent as `normalize.py` if it represents the normalization stage. Each pipeline stage module must use a lowercase snake\_case name that reflects the stage (e.g. “extract”, “transform”, “normalize”, “validate”, “write”)<sup>2</sup>.

All the above naming issues should be addressed to satisfy the naming policy. This will ensure consistency across documentation and code, and will allow automated naming lint checks (rule **ENF-001**) to pass<sup>12</sup> <sup>13</sup>.

## Duplicate Code & Architectural Patterns

There are a few areas where pipeline code and structure exhibit duplication or repeated patterns. Identifying these can help improve maintainability by refactoring into common utilities or base classes (instead of introducing new pipelines or entities<sup>10</sup>):

- **Missing ID checks in validation:** Multiple pipelines need to ensure that key identifiers are present and not null before final output. For example, the ChEMBL **Assay** pipeline’s overridden `validate` method explicitly checks for missing `assay_chembl_id` and raises a `SchemaError` if any are null<sup>14</sup>. This is a pattern likely required in other pipelines (e.g., ensuring `activity_id` or other primary keys are not missing in Activity data, ensuring parent IDs in `TestItem`, etc.). Instead of each pipeline implementing a similar null-ID check, this rule can be abstracted into the common pipeline base. **Recommendation:** provide a common validation mixin or a base class method that checks required ID fields (perhaps driven by a class attribute like `required_sort_fields` already present<sup>15</sup>). This way, all pipelines automatically enforce non-null IDs without duplicating code. For example, `ChEMBLCommonPipeline` could implement a generic check in its `validate()` and pipelines would simply declare which field(s) are critical. This reduces duplicated validation logic and ensures consistency in error handling across pipelines.
- **Output writing and artifact handling:** Many pipelines have similar logic for writing outputs and planning artifact file names. The Activity pipeline defines a custom `ActivityWriteService` and `ActivityArtifactPlanner` to handle output file naming<sup>16</sup><sup>17</sup>, and the Assay pipeline overrides `save_results` to attempt using a `PipelineOutputService` with a fallback to the base method<sup>18</sup>. These patterns – ensuring deterministic output paths and unified writing – appear in multiple places. They should be centralized. The project already has a `UnifiedOutputWriter` and a `PipelineOutputService`<sup>19</sup> for general use; however, the duplication suggests some pipelines bypass or extend them similarly. **Recommendation:** consolidate the output saving strategy in the core `PipelineOutputService` or the base pipeline so that individual pipelines do not each implement their own try/except around saving. If certain pipelines need custom file naming (as Activity does with an `activity_...csv` prefix<sup>20</sup>), provide a hook or parameter in the base planner/writer to inject name prefixes or patterns. By moving this logic to a shared service or base class, we eliminate repeated code and ensure all pipelines produce outputs in a consistent, deterministic way (satisfying idempotence and naming conventions for artifacts).

- **Normalization and parsing routines:** The ETL pipelines often perform similar normalization or parsing of nested data structures. For example, the Assay pipeline has an `AssayNormalizer` to clean and format assay fields <sup>21</sup> and an `AssayPayloadParser` to parse nested JSON payloads <sup>22</sup>. Other pipelines (Activity, Document, Target, etc.) may have analogous tasks – e.g., normalizing common fields (like chemical identifiers, organism names) or parsing nested JSON responses. If multiple pipelines implement similar normalization steps (such as flattening nested parameters or filling missing fields with placeholders like `"qc:missing_*`" as seen in Assay's `_ensure_target_integrity` <sup>23</sup>), these should be abstracted. **Recommendation:** introduce shared helper functions or a utility module in `bioetl/core` for frequently used normalization operations (e.g., a function to flatten nested JSON fields or to fill missing foreign keys with a default tag). Alternatively, a base **Normalizer** class could be defined (perhaps building on the “Base ChEMBL Normalizer” documented in common components <sup>24</sup>) so that specific pipeline normalizers can inherit standard behavior and only override what's unique. This would remove duplicate code across `normalizers.py` in different pipelines and ensure consistency (e.g., all pipelines use the same approach to type conversion, missing value handling, etc.).
- **External data enrichment pattern:** Some pipelines augment data by calling external providers (e.g., the **TestItem** pipeline fetches additional compound info from PubChem <sup>25</sup> <sup>26</sup>, and the **Target** pipeline enriches targets with UniProt/IUPHAR data as noted in its overview). The mechanism to integrate a secondary API in the transform stage is conceptually similar across these pipelines. Currently, each pipeline likely implements this on its own (TestItem's `transform()` calls PubChem, Target's pipeline presumably calls UniProt/IUPHAR in a custom step). **Recommendation:** define a generalized “enrichment” interface or mixin that pipelines can use for external data merges. For instance, a base class or utility could handle batching requests to external APIs, caching results, and merging them into the primary DataFrame. Each pipeline would then only supply the specifics (which external client to use, which fields to join on). This would reduce repeated patterns in how pipelines manage pagination, rate limiting, or merging external data. It's noted that the codebase already has a `ChembldescriptorFactory` and strategy classes for extraction <sup>27</sup> – a similar approach could be taken for enrichment strategies if multiple pipelines require it. This abstraction would make pipelines more uniform and reduce one-off implementations of external API calls scattered in different transform methods.
- **Schema and data model repetition:** The Pandera schema definitions for each entity (Activity, Assay, TestItem, etc.) might share common sub-structures (for example, all have some common audit fields or chemical identifiers). If the schema definitions (and their documentation) duplicate these shared parts, consider factoring them into a common schema or using Pandera's schema composition. For instance, if **AssaySchema**, **DocumentSchema**, etc., each define a field for `record_id` or timestamp in the same way, that could be moved to a base schema or a mixin schema. The documentation already references a **Schema Registry** <sup>28</sup>; ensuring that common fields are defined once in the registry (and reused in each schema) would prevent divergence. Also, if the documentation of schemas (under `docs/02-pipelines/schemas/...`) repeats content found in code, it may be beneficial to generate those parts of docs from the source or maintain a single source of truth for field definitions to avoid drift.

In summary, many of these duplications stem from each pipeline being implemented in isolation. By leveraging the `core` module (e.g., `bioetl.core` and `bioetl.pipelines.common`) for shared logic, the pipelines can remain thin and focused. Introduce common helpers for recurring tasks (validation checks, output formatting, external API integration) rather than copying those implementations. This will improve maintainability and ensure that fixes or improvements in one place benefit all pipelines.

# Pipeline Documentation Gaps and Improvements

Overall, the pipeline documentation is comprehensive in places but inconsistent in coverage. The documentation standards call for each pipeline's docs to cover all ETL stages and any special logic <sup>29</sup>. Upon review, some pipelines lack certain sections or depth. Below are the identified gaps and recommended improvements for documentation:

- **Coverage of all ETL stages:** Each pipeline should have documentation for its **overview**, **extract**, **transform**, **validate**, **normalize**, **write**, and **run** (execution) aspects. Currently, this coverage varies:
  - The **ChEMBL Activity** pipeline is well-documented with multiple files (overview, extraction, transformation, parser, normalizer, batch plan, column spec/mapping, schema, artifacts, etc.), effectively covering all stages except an explicit "validate" doc. It relies on a schema doc for validation. **Recommendation:** Add a brief validation section in the overview or a small `validate.md` note referencing the Pandera schema (so that the existence of a validation stage is explicit in the docs sequence). Aside from that, the Activity docs are thorough. Ensure the **run** step (how to execute the pipeline) is mentioned – e.g., pointing out it can be run via CLI and referencing the CLI documentation if not already done.
  - The **ChEMBL Assay** pipeline docs cover the overview and key components (methods, normalizer, payload parser) <sup>30</sup>, but lack a dedicated **Extraction** stage document. There is no `assay-chembl-extraction.md` outlining how data is pulled from ChEMBL for assays (likely because it uses a common extraction service). **Recommendation:** Include an Extraction doc or section stating that Assay uses the standard ChEMBL extraction mechanism (via descriptor/service) – even a short page that refers to the common extraction service would maintain consistency with other pipelines that have an extraction doc. Similarly, no separate **Transformation** doc exists; however, much of the transform logic might just be the normalizer and parser which are documented. It may be acceptable, but to mirror the pattern, consider an explicit transformation doc (even if it says "uses base transformation except for normalizer X and parser Y, see their docs"). This makes the documentation set complete and easier to navigate. The **validate** stage for Assay is implied via the schema reference but not described beyond that; a note in the overview about using `AssaySchema` for validation (which is partly there in the description) suffices. Finally, include a note on how to **run** the Assay pipeline if there are any specifics (if none, it can defer to the general CLI usage).
  - The **ChEMBL Target** pipeline documentation is minimal. We see an overview and a normalizer stub doc <sup>31</sup> <sup>32</sup>. Missing are details on **Extraction** (how target data is fetched – and importantly, how the pipeline integrates UniProt/IUPHAR data as mentioned in the overview), and **Transformation** (besides normalization, if any additional enrichment or merging logic exists). Also, since the normalizer is currently a placeholder throwing `NotImplementedError` <sup>32</sup>, the docs should clarify the status: the normalizer is intended but not yet implemented.  
**Recommendation:** Expand Target's documentation by adding an extraction page describing how ChEMBL target data is obtained (likely via the common pipeline base and possibly a specialized descriptor if one exists). Document the **enrichment process** – for instance, if after extraction the pipeline calls external APIs for extra target info, that process should be explained in a transform or enrichment section. If the implementation is pending, the docs can outline the planned behavior. Also, ensure the **validate** stage is mentioned (e.g., if `TargetSchema` exists, note that). Bringing these sections up to parity with Activity/Assay docs will make the Target documentation more complete.
  - The **ChEMBL Document** pipeline has only an overview and a "methods" doc <sup>33</sup>. This suggests a gap in documenting transform/enrichment and validation. The overview mentions enrichment from external sources, but there is no detailed doc on what those external sources are or how that works. **Recommendation:** Add a **Methods/Transformation** doc that describes any special

parsing or data merge steps for Document (for example, if it calls PubMed or CrossRef to enrich documents, detail that process). Also, if Document pipeline uses a schema for validation, it should be referenced (perhaps a DocumentSchema). In general, outline each stage: extraction (likely straightforward via base class), transformation (especially any external data added), and validation (Pandera schema). This will mirror the structure used in Activity's docs.

- The **ChEMBL TestItem** pipeline documentation is relatively extensive (overview, private methods, thin pipeline variant, CLI options, enrichment steps, etc. in its INDEX<sup>34</sup>). It covers most aspects: the **Overview** describes the pipeline and PubChem enrichment<sup>35 26</sup>, there's a doc for internal **Methods**, one for a variant pipeline without PubChem ("thin pipeline"), a **Pipeline Options** doc for run/CLI differences, and multiple docs on enrichment results and stats<sup>36</sup>. This is good, but it does break the naming convention as noted. One improvement could be to explicitly have a doc labeled **Extraction** (if the pipeline uses a descriptor or special extraction, though likely it's standard ChEMBL extraction) and **Transformation** (to compile the enrichment process). Currently, the enrichment is described across several docs (04–07 cover aspects of parent ID enrichment). It might help to clearly label one of them as the transform stage document. Perhaps the "Parent Enrichment Preparation" and "Result" docs collectively serve that purpose. **Recommendation:** Ensure the TestItem INDEX makes it clear which doc corresponds to which stage (for example, prefixing or labeling 04/05/06/07 under a Transform/Enrichment section, which it partially does with a section header). Also, for consistency, the **Normalize** stage for TestItem (if any separate normalization aside from PubChem data) isn't explicitly singled out – if the pipeline does any normalization of TestItem fields, consider documenting it. Otherwise, it's implied in transform. Overall, TestItem's documentation is detailed; just fix the naming and maybe reorganize the stage sections for clarity.
- **Consistency and formatting:** All pipeline docs should follow a uniform structure as per the documentation standards<sup>29</sup>. Common sections we expect are: **Описание** (Description/Purpose), **Модуль** (the code file path), **Наследование** (what it inherits from), **Основные методы** (key methods or stage overrides), possibly **Использование** (Usage examples), and **Related Components**. We found that most docs do include these sections. For instance, the Assay overview has Description, Module, Inheritance, Methods, Internal methods, Related Components<sup>37 38</sup>. The Target normalizer doc includes Description, Module, Methods, Usage, Related Components<sup>39 40</sup> – which is good. Any documents that are missing a section should be updated. For example, if an Overview doc didn't list the module or inheritance, add those for completeness. We did not see glaring format omissions, so the main consistency issue is ensuring **every pipeline document covers its scope in a similar level of detail**. If a pipeline has "nothing to say" for a given stage (e.g., no custom extract logic), it's better to state that briefly (or reference the base implementation) rather than omit it entirely. This way, the reader isn't left guessing whether the documentation is incomplete or the stage is intentionally identical to the base.
- **Documentation vs. code synchronization:** We observed that some documentation pages essentially duplicate code logic in narrative form (e.g., listing method signatures and describing their behavior step-by-step, as in Assay and TestItem overview docs<sup>41 25</sup>). While this is useful for understanding, it creates a maintenance burden – if method signatures or logic change in code, the docs must be updated in parallel. To mitigate this, consider using tools or conventions to keep them in sync. For example, critical constants (like field names, default values) that appear in docs could be sourced from a single config file, or code comments could be used to auto-generate parts of the docs. At minimum, reviewers should be vigilant when code changes to update the corresponding docs. A positive practice in the current docs is referencing the exact code module and class (e.g., pointing to `src/bioetl/pipelines/chembl/assay/run.py`) or

the class name), which helps trace back to implementation. **Recommendation:** Continue this practice and possibly expand it – for instance, if there's a table of column mappings in the doc and a dictionary in code for the same mapping, ensure only one source is authoritative (the code or the doc) and update the other accordingly. Inline comments in code could remind developers to update docs when a function's behavior changes. By reducing duplication (as discussed earlier) and clearly marking sections that are likely to evolve, the documentation can remain correct and up-to-date.

In summary, improving the pipeline documentation involves filling in missing pieces and enforcing consistency. Each pipeline's doc set should read like a complete story of that ETL process from start to finish, even if some parts delegate to common systems. By addressing the noted gaps (especially in Assay, Target, and Document pipelines) and fixing the structural naming issues, the documentation will better serve users and maintainers of the ETL system.

---

1 2 8 9 10 12 13 29 11-naming-policy.md

[https://github.com/SatoryKono/bioactivity\\_data\\_acquisition/blob/98f17f432532cddd56d0a07fd4796b37c54677ec/\\_docs/styleguide/11-naming-policy.md](https://github.com/SatoryKono/bioactivity_data_acquisition/blob/98f17f432532cddd56d0a07fd4796b37c54677ec/_docs/styleguide/11-naming-policy.md)

3 6 7 24 27 28 31 33 INDEX.md

<https://github.com/SatoryKono/BioactivityDataAcquisition/blob/92adfb86250a38fada7dc385bea69754c32e12a/docs/02-pipelines/chembl/INDEX.md>

4 30 INDEX.md

<https://github.com/SatoryKono/BioactivityDataAcquisition/blob/92adfb86250a38fada7dc385bea69754c32e12a/docs/02-pipelines/chembl/assay/INDEX.md>

5 34 36 INDEX.md

<https://github.com/SatoryKono/BioactivityDataAcquisition/blob/92adfb86250a38fada7dc385bea69754c32e12a/docs/02-pipelines/chembl/testitem/INDEX.md>

11 32 39 40 01-target-normalizer.md

<https://github.com/SatoryKono/BioactivityDataAcquisition/blob/92adfb86250a38fada7dc385bea69754c32e12a/docs/02-pipelines/chembl/target/01-target-normalizer.md>

14 15 18 23 run.py

[https://github.com/SatoryKono/bioactivity\\_data\\_acquisition/blob/98f17f432532cddd56d0a07fd4796b37c54677ec/src/bioetl/pipelines/chembl/assay/run.py](https://github.com/SatoryKono/bioactivity_data_acquisition/blob/98f17f432532cddd56d0a07fd4796b37c54677ec/src/bioetl/pipelines/chembl/assay/run.py)

16 17 19 20 run.py

[https://github.com/SatoryKono/bioactivity\\_data\\_acquisition/blob/98f17f432532cddd56d0a07fd4796b37c54677ec/src/bioetl/pipelines/chembl/activity/run.py](https://github.com/SatoryKono/bioactivity_data_acquisition/blob/98f17f432532cddd56d0a07fd4796b37c54677ec/src/bioetl/pipelines/chembl/activity/run.py)

21 02-assay-normalizer.md

<https://github.com/SatoryKono/BioactivityDataAcquisition/blob/92adfb86250a38fada7dc385bea69754c32e12a/docs/02-pipelines/chembl/assay/02-assay-normalizer.md>

22 03-assay-payload-parser.md

<https://github.com/SatoryKono/BioactivityDataAcquisition/blob/92adfb86250a38fada7dc385bea69754c32e12a/docs/02-pipelines/chembl/assay/03-assay-payload-parser.md>

25 26 35 00-testitem-chembl-overview.md

<https://github.com/SatoryKono/BioactivityDataAcquisition/blob/92adfb86250a38fada7dc385bea69754c32e12a/docs/02-pipelines/chembl/testitem/00-testitem-chembl-overview.md>

37 38 41 00-assay-chembl-overview.md

<https://github.com/SatoryKono/BioactivityDataAcquisition/blob/92adfb86250a38fada7dc385bea69754c32e12a/docs/02-pipelines/chembl/assay/00-assay-chembl-overview.md>