

Sintesi documentazione

README.md

- Pipeline genera KB Markdown con pre-onboarding, tagging e onboarding completo
- Output per cliente isolato (`raw/`, `book/`, `semantic/`, `config/`, `logs/`)
- Variabili d'ambiente sensibili: `SERVICE_ACCOUNT_FILE`, `DRIVE_ID`, `GITHUB_TOKEN`,
`LOG_REDACTION`, `ENV`, `CI`
- Log centralizzati, mascheramento segreti e path-safety (`is_safe_subpath`)
- Exit codes tipizzati (`ConfigError`, `PreviewError`, `PushError`)

docs/architecture.md

- RAW sempre locale; Drive solo nel pre-onboarding
- Separazione ruoli: orchestratori (UX/CLI) vs moduli tecnici senza prompt/exit
- Repository strutturato per moduli riusabili e adapter uniformi
- Invarianti: path-safety, scritture atomiche, log strutturati con `run_id`

docs/developer_guide.md

- Redazione log centralizzata via `compute_redact_flag`
- `ClientContext.load` come entry-point unico, con helper `_load_env` e `compute_redact_flag`
- Orchestratori solo UX; moduli tecnici senza `sys.exit()` o `input()`
- Principi: modularità, idempotenza, scritture sicure, consistenza API

docs/coding_rule.md

- Obbligo di type hints, docstring brevi, PEP 8, import ordinati
- Orchestratori vs moduli: `sys.exit()` solo negli orchestratori
- Path traversal: usare `is_safe_subpath`; scritture atomiche `safe_write_*`
- Config via `yaml.safe_load`, env centralizzate, cache invalidabile
- Comandi esterni tramite `proc_utils.run_cmd` con timeout/retry

docs/policy_push.md

- Push incrementale default; force push solo con `--force-push` + `--force-ack` e allow-list branch

docs/versioning_policy.md

- SemVer: MAJOR incompatibile, MINOR feature, PATCH bugfix; ogni PR che tocca API/CLI aggiorna CHANGELOG e doc corrispondenti

docs/user_guide.md

- Modalità interattiva vs CLI; pre-onboarding, tagging, onboarding con preview/push opzionali
- Log centralizzati, exit codes tipizzati, troubleshooting Docker/push/slug

Decisioni architetturali vincolanti

- Unico punto d'ingresso `ClientContext.load`; orchestratori separati dalla logica
- Redazione log gestita in `env_utils.compute_redact_flag` e propagata nei logger
- Path-safety e scritture atomiche obbligatorie per ogni output
- Push GitHub sicuro by default; force push governato da double-flag e allow-list

Mappatura codebase

Modulo/Pacchetto	Responsabilità principali	Dipendenze salienti
<code>src/pre_onboarding.py</code>	Setup locale e Drive, gestione UX/CLI, exit codes	<code>pipeline.context, pipeline.drive_utils, pipeline.config_utils</code>
<code>src/tag_onboarding.py</code>	Copia/Download PDF, generazione tags, validazione	<code>pipeline.drive_utils, pipeline.path_utils, pipeline.file_utils</code>
<code>src/onboarding_full.py</code>	Conversione PDF→MD, enrichment, preview, push	<code>pipeline.content_utils, adapters.preview, pipeline.github_utils</code>
<code>pipeline/context.py</code>	Dataclass <code>ClientContext, load</code> <code>env/config, path init</code>	<code>env_utils, path_utils, logging_utils</code>
<code>pipeline/path_utils.py</code>	Validazione slug, path traversal guard, sanitize	<code>yaml, logging_utils</code>
<code>pipeline/file_utils.py</code>	Scritture atomiche e guardie path (duplicati)	<code>path_utils (should), logging_utils</code>
<code>pipeline/env_utils.py</code>	Lettura variabili d'ambiente, redazione log	<code>dotenv, ConfigError</code>
<code>pipeline/content_utils.py</code>	Conversione e generazione markdown, fingerprint, summary/readme	<code>logging_utils, config_utils, path_utils</code>
<code>pipeline/github_utils.py</code>	Push GitHub con incrementale/force governance	<code>proc_utils, env_utils, path_utils</code>

adapters/preview.py	Avvio/stop preview Docker	proc_utils, logging_utils
semantic/*	Estrazione e validazione tag semantici	yaml, logging_utils

Doc ↔ Codice Matrix

Regola	Previsto	Osservato	Esito	Note
Path traversal prevention su disco	is_safe_subpath o ensure_within ovunque si scriva markdown verifica con	content_utils.generate_summary_is_safe_subpath	Parziale	usa is_safe_subpath ma non ensure_within
Scritture atomiche per config	safe_write_* per file critici per config	_ensure_config usa shutil.copy	Non aderente	manca atomicità
Orchestration	ori isolati da logica tecnica	sys.exit solo negli orchestratori	correttamente	Aderente moduli tecnici non chiamano exit
Redazione	log centralizzata con a	compute_redact_flag e logger redact_logs	tag_onboarding_main calcola se assente	Aderente
CSV	generati in modo sicuro	usare libreria csv	_emit_tags_csv scrive CSV manualeNon senza quoting	aderente rischio separatori

Tabelle di review

A1 – Correttezza tecnica

Principio (fonte)	Esempio nel codice	Raccomandazione
-------------------	--------------------	-----------------

Path-safety e atomicità obbligatorie	<code>file_utils.ensure_within</code> duplicato rispetto a <code>path_utils.ensure_within</code> (<i>Da migliorare</i>)	Consolidare <code>ensure_within</code> in <code>path_utils</code> ed eliminarne la copia in <code>file_utils</code> .
--	--	--

Config template scritto senza atomicità	<code>_ensure_config</code> usa <code>shutil.copy</code> (<i>Da migliorare</i>)	Copiare il template con <code>safe_write_bytes</code> e <code>ensure_within</code> .
---	---	---

Redazione log deterministica	<code>compute_redact_flag</code> implementato correttamente (<i>OK</i>)	—
---------------------------------	--	---

Attività suggerita Unificare `ensure_within` e applicare atomicità al config

Inizia l'attività

A2 – Uso corretto di funzioni/interfacce

Principio (fonte)	Esempio nel codice	Raccomandazione
Orchestratori gestiscono solo <code>pre_onboarding.py</code> usa <code>ensure_valid_slug</code> e <code>sys.exit</code> UX/exit	<code>_load_env</code> usa <code>get_env_var</code> e <code>get_bool</code> (<i>OK</i>)	—
Env centralizzate via <code>env_utils.get_env_var</code>	<code>_emit_tags_csv</code> scrive CSV manuale senza quoting (<i>Da migliorare</i>)	—
Output moduli batch-safe	<code>_emit_tags_csv</code> scrive CSV manuale senza quoting (<i>Da migliorare</i>)	Usare <code>csv.writer</code> con quoting minimale.

Attività suggerita Refactor `'_emit_tags_csv'` con `'csv.writer'`

Inizia l'attività

A3 – Ottimizzazione dei flussi

Principio (fonte)	Esempio nel codice	Raccomandazione
----------------------	--------------------	-----------------

Import ordinati
e al top-level

```
import shutil dentro _copy_local_pdfs_to_raw (Da migliorare)
```

Spostare l'import in testa al file.

Path-safety

esplícita per output usano solo `is_safe_subpath` (*Parziale*)

```
generate_summary_markdown e generate_readme_markdown
```

Applicare `ensure_within` sui file generati per fail-fast.

```
markdown
```

Attività suggerita Pulizia import e guardia path nella generazione markdown

Inizia l'attività

B – Eleganza & Best practice

Principio	Metica/Evidenza	Rischio	Impatto atteso	Sforzo
Duplicazione helper (<code>ensure_within</code>)	due implementazioni divergenti	Medio	Riduce bug e aumenta coerenza win	Quick
CSV manuale senza quoting	<code>_emit_tags_csv</code> stringhe concatenate	Alto	Migliora robustezza I/O	Quick win
Import dentro funzione	<code>import shutil interno</code>	Basso	Performance e stile	Quick win
Path guard mancante in generatori MD	solo <code>is_safe_subpath</code> senza <code>ensure_within</code>	Medio	Fail-fast su path traversal	Medio
Config copy non atomica	<code>shutil.copy in _ensure_config</code>	Medio	Evita config corrotte	Medio

Executive Summary

Valutazione complessiva: 82/100

Punti di forza

1. Architettura modulare con orchestratori separati dalla logica tecnica.
2. Logging strutturato e redazione centralizzata ben implementati.
3. Test coverage di moduli critici (path_utils, logging_utils, github_utils).

Rischi principali

1. Duplicazione della guardia `ensure_within` può generare incoerenze.
2. Serializzazione CSV manuale suscettibile a errori di parsing.
3. Scrittura `config.yaml` non atomica rischia file parziali.

Quick wins commit-ready

- Centralizzare `ensure_within`.
- Refactor `_emit_tags_csv` con `csv.writer`.
- Spostare import `shutil` e aggiungere `ensure_within` nei generatori markdown.

Prioritizzazione

- Alta: unificazione `ensure_within`, CSV sicuro.
- Media: atomicità `_ensure_config`, guardia path generatori MD.
- Bassa: pulizia import.

Roadmap

1. Refactor sicurezza I/O (Alta) – consolidare path-safety e atomic write.
2. Hardening flussi tagging (Media) – CSV robusto e validazioni aggiuntive.
3. Polish & tests (Bassa) – cleanup import, aumentare test per file_utils/env_utils.

Proposte operative (patch)

1. Deduplicare `ensure_within`

```
--- src/pipeline/file_utils.py +++ src/pipeline/file_utils.py @@ -from
.exceptions import ConfigError -from .logging_utils import
get_structured_logger --_logger =
get_structured_logger("pipeline.file_utils") -- -def ensure_within(base:
Path, target: Path) -> None: - """Garantisce che `target` risieda sotto
`base` una volta risolti i path.""" - try: - base_r = Path(base).resolve() -
tgt_r = Path(target).resolve() - except Exception as e: - raise
ConfigError(f"Impossibile risolvere i path: {e}", file_path=str(target)) from
e - - if not str(tgt_r).startswith(str(base_r)): - raise ConfigError(f"Path
traversal rilevato: {tgt_r} non è sotto {base_r}", file_path=str(target))
+from .exceptions import ConfigError +from .logging_utils import
get_structured_logger +from .path_utils import ensure_within # reuse
canonical implementation ++_logger =
get_structured_logger("pipeline.file_utils")
```

Riduce duplicazioni e mantiene una sola fonte di verità per la validazione path.

Attività suggerita Convergere su `ensure_within` in `file_utils`

Inizia l'attività

2. Scrrittura config atomica

```
--- src/pipeline/context.py +++ src/pipeline/context.py @@ -
config_path.parent.mkdir(parents=True, exist_ok=True) - template_config =
Path("config") / "config.yaml" - if not template_config.exists(): +
config_path.parent.mkdir(parents=True, exist_ok=True) + template_config =
Path("config") / "config.yaml" + if not template_config.exists(): raise
ConfigError( f"Template config.yaml globale non trovato: {template_config}",
slug=slug, file_path=template_config, ) - shutil.copy(template_config,
config_path) + from pipeline.file_utils import safe_write_bytes +
safe_write_bytes(config_path, template_config.read_bytes(), atomic=True)
return config_path
```

Garantisce scrittura atomica del config.yaml iniziale.

Attività suggerita Scrittura atomica di config.yaml

Inizia l'attività

3. CSV robusto nel tagging

```
--- src/tag_onboarding.py +++ src/tag_onboarding.py @@ - lines =
["relative_path,suggested_tags"] - for rel, tags in rows: -
lines.append(f"{rel},{tags}") - safe_write_text(csv_path, "\n".join(lines) +
"\n", encoding="utf-8", atomic=True) + import csv + with csv_path.open("w",
encoding="utf-8", newline="") as f: + writer = csv.writer(f) +
writer.writerow(["relative_path", "suggested_tags"]) + for rel, tags in rows:
+ writer.writerow([rel, tags])
```

Evita problemi con virgole nei percorsi.

Attività suggerita CSV con quoting standard

Inizia l'attività

4. Import al top-level e guardia path MD

```
--- src/tag_onboarding.py +++ src/tag_onboarding.py @@ -from pathlib import
Path -from typing import Optional, List +from pathlib import Path +from
typing import Optional, List +import shutil @@ - import shutil for src in
pdfs: try: rel = src.relative_to(src_dir)

--- src/pipeline/content_utils.py +++ src/pipeline/content_utils.py @@
summary_path = md_dir / "SUMMARY.md" - if not is_safe_subpath(summary_path,
md_dir): - raise PipelineError( - f"Path di output non sicuro per SUMMARY.md:
{summary_path}", - slug=context.slug, - file_path=summary_path, - ) +
ensure_within(md_dir, summary_path) @@ readme_path = md_dir / "README.md" -
if not is_safe_subpath(readme_path, md_dir): - raise PipelineError( - f"Path
```

```
di output non sicuro per README.md: {readme_path}", - slug=context.slug, -  
file_path=readme_path, - ) + ensure_within(md_dir, readme_path)
```

Mantiene import ordinati e applica fail-fast sulla generazione markdown.

Attività suggerita Pulizia import e path guard

Inizia l'attività

5. Test file_utils assenti → aggiungerli

```
--- tests/test_file_utils.py (nuovo) +++ tests/test_file_utils.py @@ +from  
pathlib import Path +import pytest + +from pipeline.file_utils import  
safe_write_text +from pipeline.path_utils import is_safe_subpath + +def  
test_safe_write_text_create_backup(tmp_path: Path): + target = tmp_path /  
"f.txt" + safe_write_text(target, "a") + safe_write_text(target, "b") + bak =  
target.with_suffix(".txt.bak") + assert bak.exists() + assert  
is_safe_subpath(bak, tmp_path)
```

Incrementa la copertura su operazioni critiche di I/O.

Attività suggerita Copertura minima per file_utils

Inizia l'attività

Checklists finali

Pre-commit

- ruff --fix
- black .
- pytest -q

PR

- Descrizione chiara delle modifiche e motivazioni tecniche
- Log/output allegati per `_emit_tags_csv` e scrittura `config.yaml`
- Casi limite: percorsi con virgolette, path traversal, force push non autorizzato
- Piano di rollback: revert commit in caso di regressioni

Prossimi passi consigliati

1. Core – Consolidare helper di path-safety e refactor tagging CSV (responsabile: maintainer pipeline).
2. Infra – Aggiungere test CI per `file_utils` e scenari di scrittura atomica.
3. QA – Integrare test di path traversal e quoting CSV nelle suite automatiche.