

PROJET STATION TV



Dossier de Conception Détaillé

SUIVI DES MODIFICATIONS DU DOCUMENT			
Version	Date	Validation	Commentaire
1.0	19/11/2025	Frédéric Chauvin	Rédaction initiale

REDACTEUR : Dorian BRISSON	VALIDATEUR : Frédéric CHAUVIN
CLIENT : Polytech Tours	ENCADRANT : Mathieu DELALANDRE

Table des mati res

PROJET STATION TV	1
OBJET DU DOCUMENT	4
<i>Objectif du document</i>	<i>4</i>
<i>Position dans la d�marche de d�veloppement</i>	<i>4</i>
<i>P�rim�tre de la conception d�taill�e</i>	<i>5</i>
ARCHITECTURE LOGICIELLE	5
<i>Structure globale du projet</i>	<i>5</i>
<i>Diagramme de classes</i>	<i>6</i>
<i>Description des packages</i>	<i>7</i>
<i>Sch�ma logique de l'architecture</i>	<i>9</i>
DESCRIPTION DES MODULES	9
<i>Module Preprocessing</i>	<i>9</i>
<i>Module Core Transcription</i>	<i>10</i>
<i>Module QoS</i>	<i>11</i>
<i>Module Export</i>	<i>11</i>
ALGORITHMES INTERNES	11
<i>Vue g�n�rale</i>	<i>11</i>
<i>Algorithme d'�quilibrage de charge</i>	<i>12</i>
<i>Gestion de l'affinit� CPU</i>	<i>13</i>
<i>Gestion m�moire adaptative</i>	<i>13</i>
<i>Transcription d'un fichier audio</i>	<i>14</i>
<i>Monitoring QoS en temps r�el</i>	<i>15</i>
<i>Calcul du throughput</i>	<i>15</i>
<i>Calcul du Word Error Rate (WER)</i>	<i>16</i>
<i>Contexte</i>	<i>16</i>
<i>Export incr�mental s�curis�</i>	<i>16</i>
<i>Strat�gie de tol�rance aux erreurs</i>	<i>17</i>
GESTION DES DONN�ES ET FORMATS	17
<i>Vue g�n�rale</i>	<i>17</i>
<i>Formats d'entr�e</i>	<i>18</i>
<i>Donn�es interm�diaires</i>	<i>18</i>
<i>Formats de sortie</i>	<i>19</i>
<i>Mod�les de donn�es internes</i>	<i>21</i>
<i>Donn�es de support</i>	<i>21</i>
<i>S�curit� des donn�es</i>	<i>22</i>
<i>Conformit� et tra�abilit�</i>	<i>23</i>
TESTS UNITAIRES	23
<i>Objectifs des tests unitaires</i>	<i>23</i>
<i>Environnement de test</i>	<i>24</i>
<i>Strat�gie de test</i>	<i>24</i>
<i>Modules test�s</i>	<i>25</i>
<i>Crit�res de r�ussite</i>	<i>27</i>
<i>Gestion des erreurs en test</i>	<i>27</i>
<i>Processus d'ex�cution</i>	<i>27</i>
<i>Reporting</i>	<i>28</i>
<i>Int�gration avec pipeline Agile</i>	<i>28</i>
GESTION DES ERREURS ET SECURITE	28
<i>Objectifs g�n�raux</i>	<i>28</i>
<i>Sources de risques identifi�es</i>	<i>29</i>
<i>Strat�gie g�n�rale de tol�rance aux erreurs</i>	<i>30</i>
<i>Isolation m�moire et s�curit� RAM</i>	<i>31</i>

<i>Sécurité des données</i>	31
<i>Gestion des erreurs système</i>	32
<i>Exemple</i>	32
<i>Logging et auditabilité</i>	32
<i>Sécurité opérationnelle</i>	33
<i>Exemple de protection :</i>	33
<i>Résilience du pipeline aux interruptions</i>	33
<i>Robustesse globale</i>	34
DEPLOIEMENT ET AUTOMATISATION.....	34
<i>Objectifs du déploiement</i>	34
<i>Prérequis matériels</i>	34
<i>Prérequis logiciels</i>	35
<i>Installation</i>	35
<i>Configuration</i>	36
<i>Automatisation de l'exécution</i>	36
<i>Automatisation des tests</i>	37
<i>Automatisation du monitoring</i>	37
<i>Gestion des erreurs automatisée</i>	37
<i>Modes d'exécution</i>	38
<i>Optimisation du déploiement</i>	38
<i>Sécurité opérationnelle</i>	39
<i>Reproductibilité</i>	39
<i>Étapes complètes de mise en production</i>	39

OBJET DU DOCUMENT

Objectif du document

Le présent document constitue le Dossier de Conception Détaillé (DCD) du projet Station TV – Système de transcription audio haute performance.

Son objectif est de fournir une description technique approfondie des éléments suivants :

- La **structure logicielle** du projet,
- Les **modules de traitement** et leurs responsabilités,
- Les **interfaces internes** et flux de données,
- Les **algorithmes d'orchestration et de parallélisation**,
- Les **formats structurés** utilisés (JSON, CSV, SRT),
- Les **mécanismes de monitoring et QoS**,
- Les **tests unitaires et d'intégration**.

Ce document sert de référence :

- Pour le développement, la maintenance et la réutilisation du système,
- Pour l'audit technique du projet,
- Pour l'extension future des fonctionnalités.

Position dans la démarche de développement

Le DCD intervient après la rédaction du DCG et la validation des documents de cadrage. Il a pour rôle de décrire précisément la mise en œuvre des composants logiciels, en cohérence avec les objectifs fonctionnels et non fonctionnels.

Il représente le niveau de spécification nécessaire à l'implémentation, et constitue un support pour :

- Les sprints de développement,
- La revue de code,
- La validation technique,
- Le transfert de connaissances.

P rim tre de la conception d taill e

Ce document couvre :

- Le code Python du pipeline Station TV,
- Les modules Preprocessing, Core, Export, QoS, Utilities,
- Les structures de donn es,
- Les algorithmes internes,
- Les interfaces entre modules.

Sont exclus :

- Les analyses fonctionnelles g n rales (d crites dans le DCG),
- Les sp cifications mat rielles,
- Les r sultats de performance et leur analyse d taill e.

ARCHITECTURE LOGICIELLE

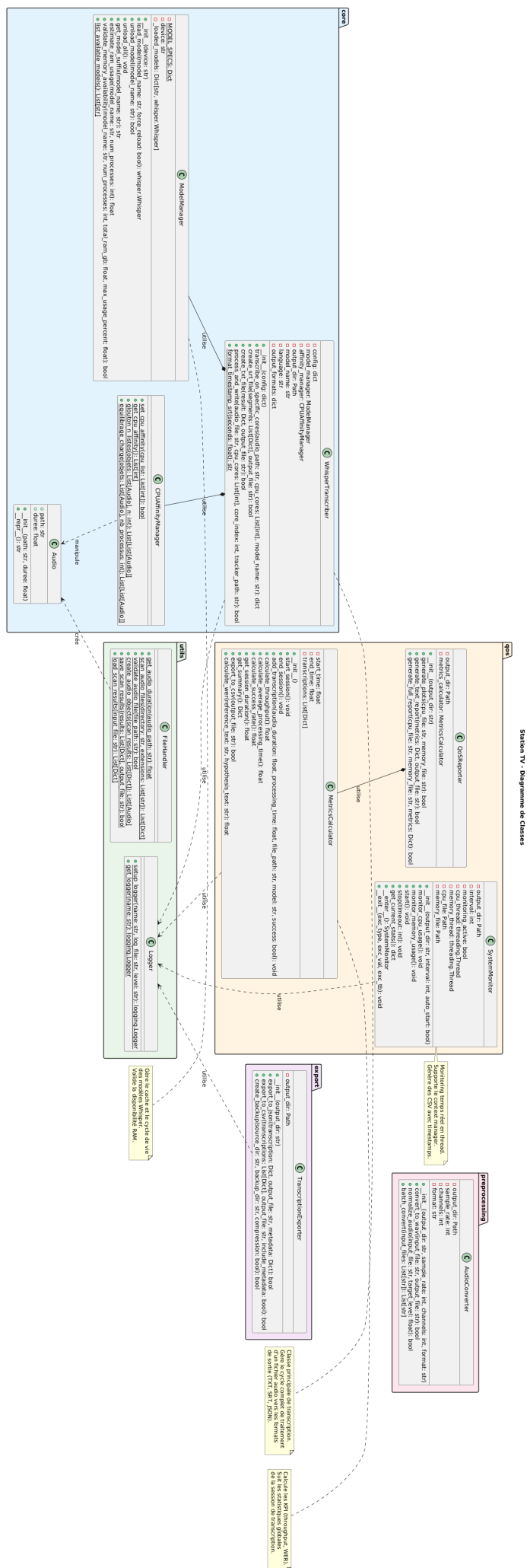
Structure globale du projet

Le code source est organis  selon une structure modulaire :

```
stationtv/  
├─ config/  
├─ core/  
├─ qos/  
├─ preprocessing/  
├─ export/  
├─ utils/  
├─ scripts/  
├─ tests/  
└─ requirements.txt
```

Chaque r pertoire correspond   un module fonctionnel ind pendant.

Diagramme de classes



Description des packages

Package preprocessing/

Contient les fonctions de préparation audio :

- Conversion MP3 → WAV
- Normalisation du signal
- Segmentation optionnelle
- Extraction de métadonnées

Fichiers principaux :

- audio_convert.py

Package core/

Implémente la transcription et l'orchestration multiprocessing.

Responsabilités :

- Gestion des modèles Whisper,
- Parallélisation CPU,
- Équilibrage de charge,
- Traitement des segments,
- Émission des métadonnées.

Fichiers :

- transcription.py
- models.py
- affinity.py

Package qos/

Responsable du monitoring et de l'analyse de performance.

Fonctionnalités :

- Mesure CPU/RAM périodique,
- Calcul throughput,
- Calcul WER,
- Export CSV + graphiques PNG.

Fichiers :

- monitor.py

- metrics.py
- reporter.py

Package export/

G re la g n ration des sorties finales :

- TXT / JSON / CSV / SRT
- sauvegardes automatiques
- int gration dans Station TV

Fichier :

- exporter.py

Package utils/

Contient les utilitaires transverses :

- logging,
- gestion fichiers et dossiers,
- chargement de configuration,
- parsing CSV.

Fichiers :

- logger.py
- file_handler.py

Package scripts/

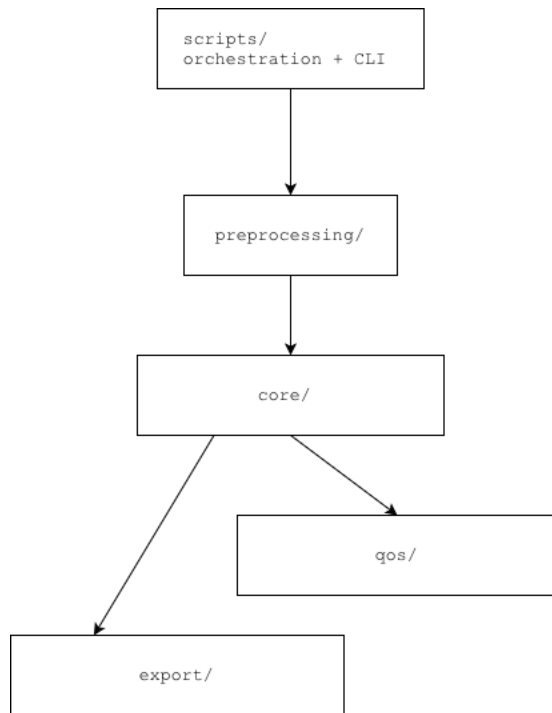
Scripts d'orchestration :

- ex cution pipeline complet,
- tests unitaires,
- batch processing.

Fichiers :

- RunPipeline.py
- RunTests.py

Schéma logique de l'architecture



DESCRIPTION DES MODULES

Module Preprocessing

Responsabilités

- standardiser les fichiers audio
- réduire la complexité acoustique
- garantir compatibilité Whisper

Entrées / Sorties

Entrée	Sortie
MP3	WAV normalisé

Fonctions principales

```
def convert_mp3_to_wav(path_in, path_out):
```

```
def normalize_audio(wav_file):
```

```
def segment_audio(wav_file, duration=15):
```

Exceptions

- invalide audio format
- erreur FFmpeg
- dur   non calculable

Module Core Transcription

Responsabilit  s

- G  rer les mod  les
- Orchestrer les processus
- G  n  rer segments + timestamps
-   viter la saturation RAM

Fonctions majeures

```
def load_model(name):
```

```
def transcribe_file(path):
```

```
def run_batch(pool, files):
```

Algorithme multi-process (simplifi  )

- 1: d  terminer N processus
- 2: mesurer RAM par mod  le
- 3: r  partir fichiers par dur  e
- 4: lancer pool
- 5: surveiller erreurs

Module QoS

Métriques

- CPU usage
- RAM usage
- throughput
- vitesse moyenne
- WER

Exports

- CSV brut
- PNG graphiques
- rapport texte

Module Export

Formats générés

Format	Usage
TXT	lecture
JSON	structuration
CSV	analyse
SRT	sous-titres

Algorithmes internes

Vue générale

Les algorithmes internes du système Station TV répondent à quatre objectifs principaux :

1. **Maximiser le débit de transcription** sur CPU multicœur.
2. **Limiter la saturation mémoire** lors du chargement de modèles lourds.
3. **Assurer une stabilité forte** lors d'exécutions longues.
4. **Produire des données exploitables** avec traçabilité complète.

Pour répondre à ces contraintes, le système combine :

- du **multi-process CPU isolé**,
- une **surveillance périodique des ressources**,
- des **exports incrémentaux sécurisés**.

Algorithme d' quilibrage de charge

Objectif

R partir un ensemble de fichiers audio entre N processus, en  quilibrant les dur es totales de traitement.

Probl me :

- Dur e inconnue a priori, mais approximable par longueur audio.
- Mod les Whisper non-streaming → m moire proportionnelle aux threads.

Approche adopt e

Par tri d croissant :

1. Trier les fichiers par dur e d croissante.
2. Cr er N listes (une par processus).
3.   chaque it ration :
 - affecter le fichier le plus long au processus le moins charg .

Pseudo-code

```
files = sort_by_duration_desc(files)
processors = [[] for _ in range(N)]
load = [0]*N
```

```
for f in files:
    idx = argmin(load)
    processors[idx].append(f)
    load[idx] += duration(f)
```

R sultat

- Variance du temps de traitement minimis e
- Charge homog ne

Gestion de l’affinit  CPU

Objectif

Allouer des c urs d di s   chaque processus pour :

- limiter les conflits,
-  viter le thrashing cache,
- assurer une stabilit  thermique et performance.

Approche

  chaque processus P :

1. D terminer un sous-ensemble de c urs disponibles.
2. D finir les masques d’affinit .
3. Appliquer via API OS.

Pseudo-code

```
cores = detect_logical_cores()
```

```
for p in processes:  
    bind(p, cores[p.index])  
    start(p)
```

Effet

- Isolation CPU par processus
- R duction significative des fluctuations

Gestion m moire adaptative

Contrainte

Le mod le Whisper Medium n cessite **>128Go RAM par instance**.
Whisper Large peut d passer **>200 Go**.

Strat gie

1. Estimer RAM disponible.
2. Calculer RAM n cessaire par mod le.
3. D terminer nombre maximal de processus.
4. Ajuster automatiquement N.

Pseudo-code

```
ram_total = detect_ram()
ram_needed = estimate_model_ram(model)
max_proc = floor(ram_total / ram_needed)

N = min(user_requested, max_proc)
```

Effet

- Crash impossibles par saturation,
- Auto-adaptation   l'environnement,
- Optimisation du d bit sans risque.

Transcription d'un fichier audio

Objectif

Convertir un fichier WAV en segments textuels avec :

- timestamps,
- scores de confiance,
- format structur .

Pseudo-code

```
model = load_model()
audio = load_audio(path)

segments = model.transcribe(audio,
                             timestamps=True,
                             word_timestamps=True)

for seg in segments:
    save(seg)
```

Structure segment

```
{
  "start": 12.96,
  "end": 15.78,
  "text": "Bonjour   tous",
  "confidence": 0.92
}
```

Monitoring QoS en temps r el

Objectif

Mesurer les m triques syst me **toutes les 30 secondes**.

M triques enregistr es

- CPU %
- RAM %
- RAM absolue
- d bit (x temps r el)
- vitesse moyenne
- WER

Pseudo-code

```
while running:  
  cpu = psutil.cpu_percent()  
  ram = psutil.virtual_memory().percent  
  save_metrics(cpu, ram)  
  sleep(30)
```

Effet

- S rie temporelle coh rente
- Diagnostic a posteriori
- Alertes pilot es par donn es

Calcul du throughput

D finition

throughput = (dur e audio trait e) / (dur e r elle)

Objectif

- throughput $\geq 1.4\times$ sur mod le Medium

Ex cution

```
start = now()  
audio_processed = 0
```

```
for f in files:  
  audio_processed += duration(f)
```

```
elapsed = now() - start  
throughput = audio_processed / elapsed
```

Calcul du Word Error Rate (WER)

Contexte

WER = indicateur de qualit  de transcription.
Mesur  sur un ** chantillon**.

Formule

$$\text{WER} = (S + D + I) / N$$

Avec :

- S = substitutions
- D = suppressions
- I = insertions
- N = nombre total de mots

Pseudo-code

```
def wer(ref, hyp):  
    S, D, I = align(ref, hyp)  
    return (S + D + I) / len(ref)
```

Export incr mental s curis 

Objectif

 viter perte de donn es en cas de crash.

M canisme

- sauvegarde par fichier trait ,
- flush syst matique,
- backups p riodiques.

Pseudo-code

```
for file in processed:  
    json_save(file)  
    txt_save(file)  
    if backup_enabled:  
        zip(file)
```


Stratégie de tolérance aux erreurs

Types d'erreurs

- audio corrompu
- crash modèle
- conflit mémoire
- I/O

Stratégie

- try/catch individuel,
- logs détaillés,
- retry automatique (max 3).

Pseudo-code

```
for f in files:
    for attempt in range(3):
        try:
            process(f)
            break
        except:
            log("error")
    else:
        log("failed permanently")
```

Gestion des données et formats

Vue générale

Le système Station TV manipule quatre catégories principales de données :

1. **Données audio source** : fichiers MP3 issus des chaînes TNT.
2. **Données intermédiaires** : WAV normalisés, segments internes, métriques QoS.
3. **Données finales** : transcriptions et exports multi-formats.
4. **Données de support** : logs, configurations, rapports, backups.

La gestion des données est organisée selon trois principes clés :

- **Standardisation des formats** pour compatibilité Whisper.
- **Structuration riche** pour exploitation aval.
- **Sécurité et traçabilité** pour opérations longues et critiques.

L'ensemble des données est traité de manière **incrémentale** pour limiter les risques de perte en cas de crash.

Formats d'entrée

Fichiers audio MP3

Propriété	Valeur
Encodage	MP3
Fréquence	48 kHz
Débit	256 kbps
Mono	Variable (mono/dual)

Métadonnées extraites automatiquement :

- nom du fichier
- durée audio
- chaîne
- horaire diffusion

Structure CSV d'inventaire générée :

```
filename,channel,start_time,duration_s,path
tf1_03112025_2030.mp3,TF1,2025-11-03 20:30,312,/data/input/tf1_0311.mp3
```

Données intermédiaires

Fichiers WAV normalisés

Propriété	Valeur
Format	PCM16
Fréquence	48 kHz
Canaux	Mono
Normalisation	Oui

Nom standardisé :

<original>_norm.wav

Segments internes JSON

Les segments transcrits par Whisper avant export sont stock s en JSON structur .

Structure interne d'un segment :

```
{  
  "start": 3.02,  
  "end": 5.17,  
  "text": "Bonsoir et bienvenue   tous",  
  "confidence": 0.87,  
  "speaker": null  
}
```

Historique des m triques QoS

Fichiers CSV par run, append temps r el :

```
timestamp,cpu,ram,throughput,wer  
2025-11-21 12:30:00,74.5,62.3,1.45,0.09
```

Formats de sortie

Le syst me produit quatre types de donn es finales.
Chaque export correspond   un use-case m tier diff rent.

TXT (texte brut)

Objectif : lecture humaine, traitement linguistique simple.

Format :

```
[00:00:02] Bonsoir   tous et bienvenue.  
[00:00:05] Ce soir dans notre  mission...
```

Caract ristiques :

- UTF-8
- timestamps optionnels
- un fichier par source audio

JSON (structure riche)

Objectif : exploitation machine, APIs, intégration Station TV.

Structure générale :

```
{  
  "file": "tf1_03112025_2030.mp3",  
  "channel": "TF1",  
  "duration": 312,  
  "model": "medium",  
  "segments": [...]  
}
```

CSV (tableau consolidé multi-sources)

Objectif : analyse statistique, import Excel, data warehouse.

Exemple :

```
file,channel,start,end,text,confidence  
tf1_x,TF1,3.02,5.17,"Bonsoir à tous...",0.87
```

Avantages :

- exploitation rapide
- agrégation simple
- format universel

SRT (sous-titres horodatés)

Objectif : vidéo, accessibilité, médias.

Exemple :

```
1  
00:00:02,000 --> 00:00:05,170  
Bonsoir à tous et bienvenue.
```

```
2  
00:00:05,170 --> 00:00:07,800  
Ce soir dans notre émission...
```

Contraintes :

- segment < 4 s recommandé
- texte lisible et fluide

Mod les de donn es internes

Repr sentation d'un fichier audio

```
class AudioFile:  
    path: str  
    duration: float  
    channel: str  
    timestamp: datetime
```

Repr sentation d'un segment transcrit

```
class Segment:  
    start: float  
    end: float  
    text: str  
    confidence: float  
    speaker: Optional[str]
```

Donn es de support

Fichier de configuration YAML

Structure typique :

```
hardware:  
    cpu_processes: 6
```

```
whisper:  
    model: medium  
    timestamps: true
```

```
export:  
    formats:  
        - txt  
        - json  
        - csv  
        - srt
```

```
logging:  
    level: INFO
```

Objectifs :

- centraliser les param tres
- permettre le tuning rapide
- assurer la reproductibilit 

Journaux d'ex cution

Caract ristiques :

- logs horodat s
- par module
- multi-niveaux : DEBUG, INFO, ERROR

Exemple :

[2025-11-21 12:42:01] INFO: Started process #2 with 64000 MB RAM

Backups

Strat gies :

- zip par fichier trait 
- flush imm diat
- export sur NAS optionnel

Nom standardis  :

backup_2025-11-21_1530.zip

S curit  des donn es

Le syst me applique des mesures de s curit  minimisant la perte d'information :

1. **Export incr mental**
Un fichier export  = un fichier sauvegard .
2. **Redondance optionnelle NAS**
3. **Protection contre crash**
 - try/catch par fichier
 - logs erreurs
 - skip non bloquant
4. **Nom standardis  + m tadonn es**
5. **Partitionnement par run**

Conformit  et tra abilit 

Chaque transcription est associ e   des m tadonn es permettant :

- la reconstitution du contexte
- la comparaison entre runs
- l’audit scientifique

Exemple JSON :

```
{  
  "file": "france2_x",  
  "date": "2025-11-21",  
  "model": "medium",  
  "ram": 6824,  
  "cpu_avg": 74.12,  
  "wer": 0.086  
}
```

Tests unitaires

Objectifs des tests unitaires

Les tests unitaires du projet Station TV ont pour objectif de :

1. V rifier le bon fonctionnement des modules critiques (Preprocessing, Core, QoS, Export, Utilities)
2. D tecter pr cocement les r gressions
3. Garantir la stabilit  du syst me dans un contexte de d veloppement Agile
4. Fournir une base solide pour l’int gration continue
5. S curiser la mont e en charge ult rieure

La priorit  est donn e :

- aux composants   forte complexit  algorithmique,
- aux fonctions manipulant des donn es critiques,
- aux fonctions expos es aux erreurs d’E/S.

Environnement de test

Langage et Framework

Technologie	Version
Python	3.10+
pytest	≥ 9.0

Organisation

Les tests sont regroup s dans le r pertoire :

tests/

Le lancement se fait via :

python scripts/RunTests.py

ou directement :

pytest -v

Strat gie de test

La strat gie adopt e repose sur trois niveaux :

Tests unitaires

- Portent sur fonctions ou classes isol es
- Validation comportemental / structurel

Tests fonctionnels

- Portent sur modules complets
- V rification coh rence I/O

Tests de non-r gression

- Assurent la stabilit  suite   modifications
- Tests d clench s avant merge Git

Modules test s

Module Preprocessing

Fonctionnalit s test es :

- conversion MP3 → WAV
- extraction de m tadonn es
- segmentation optionnelle
- traitement fichiers invalides

Exemples de tests :

```
def test_convert_mp3_to_wav(tmp_path):  
    wav = convert_mp3_to_wav("test.mp3", tmp_path/"test.wav")  
    assert wav.exists()
```

Crit res :

- dur e correcte
- absence d'erreur

Module Core Transcription

Fonctionnalit s test es :

- chargement de mod le Whisper
- transcription courte
- gestion multiprocessing
- erreur fichier corrompu

Exemples :

```
def test_transcribe_short_clip():  
    text = transcribe_file("short.wav")  
    assert len(text) > 0
```

Crit res :

- pas de crash
- sortie non vide
- timestamps ordonn s

Module QoS

Fonctionnalit s test es :

- mesure CPU/RAM
- calcul throughput
- calcul WER
- export CSV

Exemple :

```
def test_compute_throughput():  
    result = compute_throughput(600, 300)  
    assert result == 2.0
```

Module Export

Fonctionnalit s test es :

- g n ration TXT
- g n ration JSON
- coh rence des timestamps
- formats valides

Crit res :

- export lisible
- absence de caract res invalides
- JSON parsable

Module Utilities

Fonctions test es :

- logging multi-niveaux
- parsing CSV
- lecture config YAML

Exemple :

```
def test_load_config():  
    config = load_config("default.yaml")  
    assert config["hardware"]["cpu_processes"] > 0
```

Crit res de r ussite

Un test unitaire est consid r  comme r ussi si :

- la fonction retourne un r sultat attendu,
- aucun crash n'est observ ,
- les donn es g n r es sont valides.

Un run complet est consid r  valide si :

1. Tous les tests passent, ou
2. Les  checs sont document s et justifi s.

Gestion des erreurs en test

Strat gies :

- tests sp cifiques pour cas extr mes
- fixtures pour donn es invalides
- simulations d'erreurs E/S

Exemple :

```
def test_corrupted_file():  
    with pytest.raises(Exception):  
        transcribe_file("null.wav")
```

Processus d'ex cution

Le lancement des tests se fait id alement :

- Avant chaque merge Git
- Avant une ex cution batch longue
- Apr s chaque modification critique

Script de test automatis  :

```
python scripts/RunTests.py
```

Sortie attendue :

15 passed, 0 failed

Reporting

Les r sultats des tests peuvent  tre :

- affich s en console,
- export s en fichier texte,
- archiv s dans Git.

Format type :

```
[TEST] test_convert_mp3_to_wav ..... OK  
[TEST] test_transcribe_short_clip ... OK
```

Int gration avec pipeline Agile

Chaque sprint int gre :

- objectifs de tests sp cifiques
- crit res de validation
- rapport de tests

Gestion des erreurs et s curit 

Objectifs g n raux

Le syst me Station TV est con u pour fonctionner sur de longues p riodes (plusieurs dizaines d'heures) et traiter un grand volume de donn es audio.

Dans ce contexte, la gestion des erreurs et la s curit  des donn es sont des enjeux majeurs.

Les objectifs sont :

1. **Assurer la continuit  d'ex cution**, malgr  des erreurs locales.
2. ** viter les crashes critiques**, notamment li s   la m moire ou aux E/S.
3. **Garantir l'int grit  et la tra abilit  des donn es produites**.
4. **Limiter les pertes de donn es en cas d'interruption soudaine**.
5. **Fournir suffisamment d'informations pour diagnostic et audit**.

Sources de risques identifiées

Les risques potentiels sont regroupés en cinq catégories :

Erreurs liées aux fichiers audio

- fichiers corrompus ou incomplets
- encodage non compatible
- durée non calculable

Erreurs logicielles

- exceptions Python non gérées
- dysfonctionnement librairie audio
- crash modèle Whisper

Erreurs système

- manque de RAM
- surcharge CPU
- saturation disque

Erreurs E/S

- permissions insuffisantes
- fichier temporaire inaccessible
- chemins invalides

Erreurs liées au monitoring

- métriques incohérentes
- échec génération graphique
- formats CSV invalides

Strat gie g n rale de tol rance aux erreurs

Le syst me adopte une politique de **r silience locale** :

Principe cl 

Une erreur sur un fichier ne doit pas arr ter l'ensemble du pipeline.

Chaque unit  de traitement (fichier audio) est g r e de fa on isol e, avec :

- protection par try/catch,
- logs d taill s,
- tentatives multiples,
- skip contr l .

Illustration simplifi e

```
for f in files:  
    try_three_times(process_file, f)
```

R sultat :

- un crash local   crash global
- robustesse aux donn es h t rog nes
- pipeline long s curis 

M canisme de retry automatique

Pour les op rations critiques (transcription, conversion), le syst me effectue jusqu'  **trois tentatives**.

Pseudo-code :

```
def try_three_times(func, arg):  
    for attempt in range(3):  
        try:  
            return func(arg)  
        except Exception as e:  
            log_error(e, attempt)  
    mark_as_failed(arg)
```

Avantages :

- r duction des erreurs transitoires
- meilleure stabilit  sur longues sessions

Isolation m moire et s curit  RAM

Pour  viter les crashes li s   la m moire, le syst me met en  uvre :

1. **Processus ind pendants**
 -   pas de partage de mod le
 -   allocation RAM isol e
2. **Dimensionnement dynamique**
 -   limite auto du nombre de processus
3. **Monitoring RAM p riodique**
4. **Arr t contr l  en cas de saturation**

Exemple de logique :

```
if ram_usage() > threshold:  
    pause_or_kill_process()
```

Effets :

- crash RAM  vit 
- stabilit  longue dur e

S curit  des donn es

Sauvegarde incr mentale

Chaque fichier transcrit est sauvegard  imm diatement.

Cons quences positives :

- perte maximale limit e   1 fichier
- reprise facile apr s incident

Syst mes de backup

Option de backup automatique :

- ZIP par session
- export hors dossier de travail

Objectif :

Tol rance   une corruption totale du dossier local.

Format de donn es robuste

- UTF-8
- JSON valid 
- CSV structur 
- SRT conforme   la norme

Gestion des erreurs syst me

D tection d'anomalies

Surveillance :

- CPU > X%
- RAM > Y%
- throughput trop faible
- erreurs r p t es

Actions correctives

- arr t d'un processus
- r duction du nombre de workers
-  criture d'un rapport d'erreur
- reprise partielle

Exemple

```
if ram_usage() > limit:  
    decrease_process_count()
```

Logging et auditabilit 

Les logs sont produits automatiquement, avec :

- timestamp
- niveau (INFO, WARNING, ERROR)
- contexte (module, fichier)
- message d taill 

Exemple :

```
[2025-11-21 12:32:15] ERROR: Transcription failed on file 'f2.wav', attempt 2
```

Objectifs :

- permettre diagnostic post-mortem
- analyser la performance du pipeline
- assurer tra abilit  scientifique

Sécurité opérationnelle

Le script RunPipeline intègre :

- validation de configuration avant exécution
- vérification des chemins d'entrée / sortie
- estimation RAM vs modèle sélectionné
- avertissement si risque de crash

Exemple de protection :

```
if required_ram > available_ram:  
    warn_user()  
    abort()
```

Effet :

- erreurs préventives
- réduction des mauvaises configurations

Résilience du pipeline aux interruptions

Cas typiques :

- crash processus
- arrêt machine
- reboot imprévu

Grâce à :

- sauvegarde incrémentale
- logs détaillés
- fichiers marqués comme "done"

Le système peut :

- reprendre après crash,
- éviter tout retraitement inutile.

Robustesse globale

Grâce à l'ensemble de ces mécanismes, Station TV est capable de :

- traiter des volumes massifs,
- sur de longues périodes,
- avec un taux d'échec marginal,
- et une perte d'information quasi nulle.

La résilience est une caractéristique structurelle du système, pas un patch tardif.

Déploiement et automatisation

Objectifs du déploiement

Le déploiement du système Station TV vise à fournir une installation :

1. reproductible
2. automatisée autant que possible
3. facilement configurable
4. robuste sur de longues exécutions

Le système a été conçu pour être installé sur des architectures similaires à la station Dell Precision 5820, mais reste compatible avec d'autres configurations, sous réserve d'ajustement des ressources.

Prérequis matériels

Recommandations minimales pour un usage intensif :

Composant	Spécification recommandée
CPU	16–24 threads
RAM	≥ 128 Go (Medium) / ≥ 200 Go (Large)
GPU	Aucun (CPU-only)
Stockage	SSD NVMe
OS	Ubuntu 22.04 ou Windows Server 2022

Cette configuration a été sélectionnée pour répondre aux contraintes RAM/CPU élevées du modèle Whisper.

Prérequis logiciels

Environnement Python

- Python 3.10+
- pip (package manager)
- virtualenv (optionnel)

Bibliothèques nécessaires

Installées via :

```
pip install -r requirements.txt
```

Logiciel tiers

- FFmpeg (conversion audio)

Vérification :

```
ffmpeg -version
```

Installation

Clonage du projet

```
git clone <repository-url>  
cd stationtv
```

Environnement virtuel

Linux :

```
python -m venv venv  
source venv/bin/activate
```

Windows :

```
python -m venv venv  
venv\Scripts\activate
```

Installation des dépendances

```
pip install -r requirements.txt
```

Configuration

Le comportement du syst me est pilot  par un fichier YAML, typiquement :

```
hardware:  
  cpu_processes: 6  
  
whisper:  
  model: medium  
  timestamps: true  
  
export:  
  formats: [txt, json, srt, csv]  
  
paths:  
  input: "./data/input"  
  output: "./data/output"  
  logs: "./logs"
```

Points cl s :

- nombre de processus ajust  automatiquement selon RAM
- choix du mod le (tiny → large)

Automatisation de l'ex cution

Script d'ex cution complet

L'ex cution du pipeline complet est r alis e en un clic via :

```
python scripts/RunPipeline.py
```

Fonctions assur es automatiquement :

1. Scan r pertoire audio
2. Pr traitement MP3 → WAV
3. Transcription multiprocessing
4. Export multi-format
5. Monitoring QoS
6. Logging
7. Sauvegardes

Un fichier batch/PowerShell est fourni pour Windows :

```
RUN_PIPELINE.bat
```

Automatisation des tests

Les tests peuvent être lancés via :

`python scripts/RunTests.py`

ou :

`pytest -v`

Utilisation recommandée :

- avant chaque exécution longue
- avant merge Git
- avant modification configuration

Automatisation du monitoring

Le monitoring est intégré au run complet.

Rôle de l'automate QoS :

1. collecter métriques CPU / RAM
2. Aggréger throughput / WER
3. Exporter CSV + PNG
4. Calculer moyennes / variances
5. Signaler anomalies possibles

Structure standard des exports :

```
/output/  
monitoring_cpu.csv  
monitoring_memory.csv  
cpu_usage.png  
memory_usage.png
```

Gestion des erreurs automatisée

Pendant l'exécution, le système :

- détecte les erreurs locales
- applique retry automatique
- journalise l'événement
- skip si échec persistant

Aucune intervention manuelle nécessaire pour :

- erreurs temporaires
- fichiers corrompus
- saturation locale

Modes d'exécution

Mode batch complet

Pour traiter un grand lot de fichiers :

```
python scripts/RunPipeline.py
```

Mode unitaire

```
python scripts/BasicTestWhisper.py file.mp3
```

Optimisation du déploiement

Dimensionnement automatique des ressources

Le système évalue :

- RAM totale disponible
- RAM requise par modèle
- nombre optimal de processus

et ajuste dynamiquement.

Affinité CPU

- isolation des workers
- réduction des conflits
- amélioration throughput

Monitoring long terme

Adapté à des sessions > 5j :

- sérialisation CSV
- graphiques haute résolution
- analyse post-mortem

S curit  op rationnelle

Protection avant ex cution :

- validation configuration YAML
- v rification chemins d'entr e / sortie
- estimation RAM > seuil
- avertissement utilisateur

Protection pendant ex cution :

- kill pr ventif
- pr vention crash RAM
- suppression deadlocks

Reproductibilit 

La reproductibilit  est assur e par :

- environnement virtuel
- versionnement Git
- configuration externalis e YAML
- documentation exhaustive

R sultat :

Ex cution identique sur n'importe quelle machine compatible.

 tapes compl tes de mise en production

1. Installer d pendances
2. Configurer YAML
3. Lancer tests unitaires
4. Lancer pipeline
5. Surveiller QoS
6. Exporter r sultats
7. Sauvegarde finale

Dur e typique :

- quelques heures   plusieurs jours selon volume / mod le