



M1 INFORMATIQUE AIGLE

HMIN201

M1 TER

---

## TER Software Heritage

Rapport Final

---

### Groupe BAJONIM

**Bachar RIMA,**

bachar.rima@etu.umontpellier.fr

**Joseph SABA,**

joseph.saba@etu.umontpellier.fr

**Tasnim SHAQURA,**

tasnim.shaqura@etu.umontpellier.fr

*Encadrant :*

Jessie CARBONNEL

*Responsable de l'UE :*

Mattieu LAFOURCADE

27 mai 2019

# Table des matières

<b>1</b>	<b>Introduction</b>	<b>3</b>
1.1	Description de Software Heritage . . . . .	3
1.2	Contexte du TER . . . . .	3
1.3	Plan du rapport . . . . .	3
<b>2</b>	<b>Problématique</b>	<b>4</b>
2.1	La diaspora du code source . . . . .	4
2.2	La fragilité du code source . . . . .	5
2.3	Software Heritage en tant que solution . . . . .	5
2.4	Les défis . . . . .	5
2.5	Les principes de base . . . . .	5
2.5.1	Transparence et gratuité . . . . .	5
2.5.2	Réplication compréhensive de l'entiereté du système . .	6
2.5.3	Multitude des partenaires, sans profits . . . . .	6
2.5.4	Pas de présélection . . . . .	6
2.5.5	Source Code First . . . . .	6
2.5.6	Identifiants Intrinsèque . . . . .	6
2.5.7	Informations de provenance et informations factuelles .	7
2.5.8	Minimalisme . . . . .	7
2.6	Notre travail . . . . .	7
<b>3</b>	<b>Analyse</b>	<b>8</b>
3.1	Fonctionnement de Software Heritage . . . . .	8
3.1.1	Modèle des données . . . . .	8
3.1.2	Architecture et flot des données . . . . .	8
3.1.3	L'archive . . . . .	8
3.2	Méthodologie . . . . .	9
3.3	Planning Prévisionnel . . . . .	9
<b>4</b>	<b>Conception</b>	<b>10</b>

<b>5</b>	<b>Implémentation</b>	<b>11</b>
<b>6</b>	<b>Résultats</b>	<b>12</b>
<b>7</b>	<b>Conclusion</b>	<b>13</b>
7.1	Planning final . . . . .	13
7.2	Difficultés rencontrées . . . . .	14
7.3	Perspectives . . . . .	14
7.4	Bilan et apports du TER . . . . .	14

# Chapitre 1

## Introduction

1.1 Description de Software Heritage

1.2 Contexte du TER

1.3 Plan du rapport

# Chapitre 2

## Problématique

Les logiciels sont actuellement omniprésents dans tous les aspects de notre vie quotidienne ; archiver leurs codes source paraît ainsi une tâche primordiale. À ce titre là, des plateformes ont déjà été proposées, telles que The Internet Archive et UNESCO Persist. Toutefois, ces plateformes se concentraient plutôt sur la préservation des fichiers exécutables ; allant jusqu'à offrir des émulateurs pour permettre l'exécution des logiciels présents dans leurs archives. Par comparaison, Software Heritage s'intéresse au code source des logiciels, pas à leurs exécutables. En effet, le code source d'un logiciel constitue un artefact logiciel fondamental dans le domaine des connaissances scientifiques, culturelles, et techniques. Le code source est écrit sous une forme compréhensible par les humains, et peut facilement être transformé en une forme exécutable par une machine. Le code source est muable et évolue selon les besoins. Sa préservation nous permet d'accéder à l'historique du développement d'un logiciel.

Malgré son importance dans notre vie quotidienne, il est facile de voir que nous prenons pas soin correctement du code source. Cela est dû à trois raisons principales.

### 2.1 La diaspora du code source

Le quantité de projets open sources a vu un énorme accroissement pendant les deux dernières décennies. Le code source de ces projets sont souvent développés sur des plateformes d'hébergement publiques (comme *Github* et *BitBucket*), ou sur des divers forges institutionnelles. Beaucoup d'options s'offrent aux développeurs pour distribuer leurs logiciels. La distribution peut se faire sur des plateformes comme *Github*. Elle peut se faire via des archives liées à des écosystèmes spécifiques, comme *CTAN*, qui distribue des

logiciels pour TeX. Les développeurs peuvent aussi choisir de publier leurs logiciels sur des distributions comme *Debian* et *Fedora*, ou via un gestionnaire de paquets comme *npm* et *pip*.

## 2.2 La fragilité du code source

Le code source est une entité fragile. Elle peut être facilement détruite ou perdue si elle n'est pas fréquemment sauvegardée. Les plateformes d'hébergement ne garantissent pas forcément la préservation de leurs contenus ; des grandes plateformes d'hébergement ont déjà arrêté leurs services.

## 2.3 Software Heritage en tant que solution

Software Heritage a été créé pour relever ces défis. Software Heritage vise à fournir une infrastructure qui permet la collection, l'organisation, la préservation, et l'accès à tout code source public. L'archive doit avoir la capacité d'accomplir ses objectifs pour toutes plateformes de développement et de distribution, et doit pouvoir persister les codes source sur le long terme.

### Current status et roadmap de SWH

## 2.4 Les défis

**Identifier les plateformes d'hébergement** : Les projets peuvent être hébergés sur des plateformes bien connues, comme sur des plateformes obscures. Il faut construire un catalog des plateformes. **Supporter différents protocoles** : Software Heritage doit pouvoir récupérer leurs projets et doit pouvoir maintenir les modifications faites sur ces projets. Vu qu les plateformes d'hébergements sont hétérogènes, Software Heritage essayera de promouvoir des bonnes pratiques pour la préservation. **Parcourir les historiques de développement** : Les plateformes d'hébergements supportent différents logiciels de gestion de versions qui n'ont pas les même modèles de données. Software Heritage construira un tel modèle unifiant.

## 2.5 Les principes de base

### 2.5.1 Transparence et gratuité

Pour pouvoir assurer la préservation de l'archive sur le long terme, il faut que tout les éléments formant l'archive soient open source et accessibles au

publique.

### **2.5.2 Réplication comprehensive de l'entiereté du système**

Un tel archive est soumis à différents types de risques. Ces risques étant inévitables, le système doit pouvoir les tolérer. Le système sera répliqué sur différents niveaux : différentes localisations géographiques, différents matériels de stockage , etc...

### **2.5.3 Multitude des partenaires, sans profits**

Pour atteindre ses objectifs, Software Heritage ne doit pas dépendre sur une seule entité qui cherche d'en profiter, et ne doit pas être créé pour la génération de profits. Ce projet doit apporter de la valeur au public en large, et non seulement pour les organisations qui le supporte.

### **2.5.4 Pas de présélection**

Il est impossible de savoir quel projets vont finir par être les plus importants. Il faut donc préserver tout les logiciels disponibles sans présélection, surtout que la capacité technique de faire cela est disponible.

### **2.5.5 Source Code First**

Bien qu'il est intéressant de garder le contexte du code source (comme les Wiki, l'environnement où le programme est exécuté, etc), une telle tâche nécessite une énorme quantité de ressources, surtout qu'il n'y a pas de préselection. Software Heritage se contente d'archiver les code sources ainsi que leur historique de développement capturé par les logiciels de gestion de versions. Cela permet de garder des informations importantes qu'on retrouve dans les messages de commit.

### **2.5.6 Identifiants Intrinsèque**

Les identifiants des objets stockés ne doivent pas dépendre des sources externes et doivent pouvoir être calculés à partir des objets qu'ils identifient. Ils sont étroitement liés à ces objets. Cela permet de vérifier que l'objet obtenu correspond à l'objet demandé, et permet la détection des modifications sur l'objet.

### **2.5.7 Informations de provenance et informations factuelles**

Software Heritage va stocké les informations de provenance qui décrivent le ou, le quoi, et le quand des objets dans l'archive. Ces informations vont être vérifiées et les méthodes de vérification vont être stockées aussi.

### **2.5.8 Minimalisme**

Software Heritage se contentera de construire l'infrastructure essentielle et rien de plus.

## **2.6 Notre travail**

Dans le cadre de ce projet, encadré par Jessie Carbonnel, nous avons cibler la plateforme d'hébergement Launchpad afin de collectionner les codes sources qui y sont hébergés, et les stocker dans l'archive de Software Heritage. Ainsi, les objectifs de ce TER peuvent être énumérés de la manière suivante :

- Lire/comprendre les articles et tutoriels écrits par l'équipe de Software Heritage ;
- Analyser différentes plateformes d'hébergement afin d'en cibler une ;
- Concevoir et développer un Lister pour la plateforme choisie ;
- Répliquer localement l'environnement de Software Heritage afin de tester le Lister développé ;
- Faire une Pull Request afin d'intégrer le Lister testé au dépôt de développement de Software Heritage.



# Chapitre 3

## Analyse

### 3.1 Fonctionnement de Software Heritage

#### 3.1.1 Modèle des données

Plateformes d'hébergement

Artéfacts logiciels

Informations sur la provenance des données

Structure des données

#### 3.1.2 Architecture et flot des données

Flot d'ingestion des données

Listing

Loading

Scheduling

#### 3.1.3 L'archive

Stockage des noeuds BLOB de l'archive

les IDs, les fichiers

Stockage des autres noeuds des archives

chemins, répertoires, snapshot, revisions, releases  
postgres, etc

**Stockage haché des objets**

**Mise en miroir des noeuds**

**Politique de rétention**

**Récupération automatique des objects corrompus**

## **3.2 Méthodologie**

Sourceforge sitemap, api

Launchpad api, client

Analyzing the listers (bitbucket, gitlab, github, eclipse, LIRMM, OpenHub, Assembla, GNU savannah

heritage, injection de dependances

conclusion : on adoptera une strategie pour definir un lister, loader ou autre

## **3.3 Planning Prévisionnel**

# Chapitre 4

## Conception

design de la solution proposée (diagrammes + explications)

# Chapitre 5

## Implémentation

Les Listers de Software Heritage et les codes qui les accompagnent sont écrits en Python.

les technos qu'on a utilisé

bibliothèques

Outils (e.g. XML parsers)

Launchpad client

# Chapitre 6

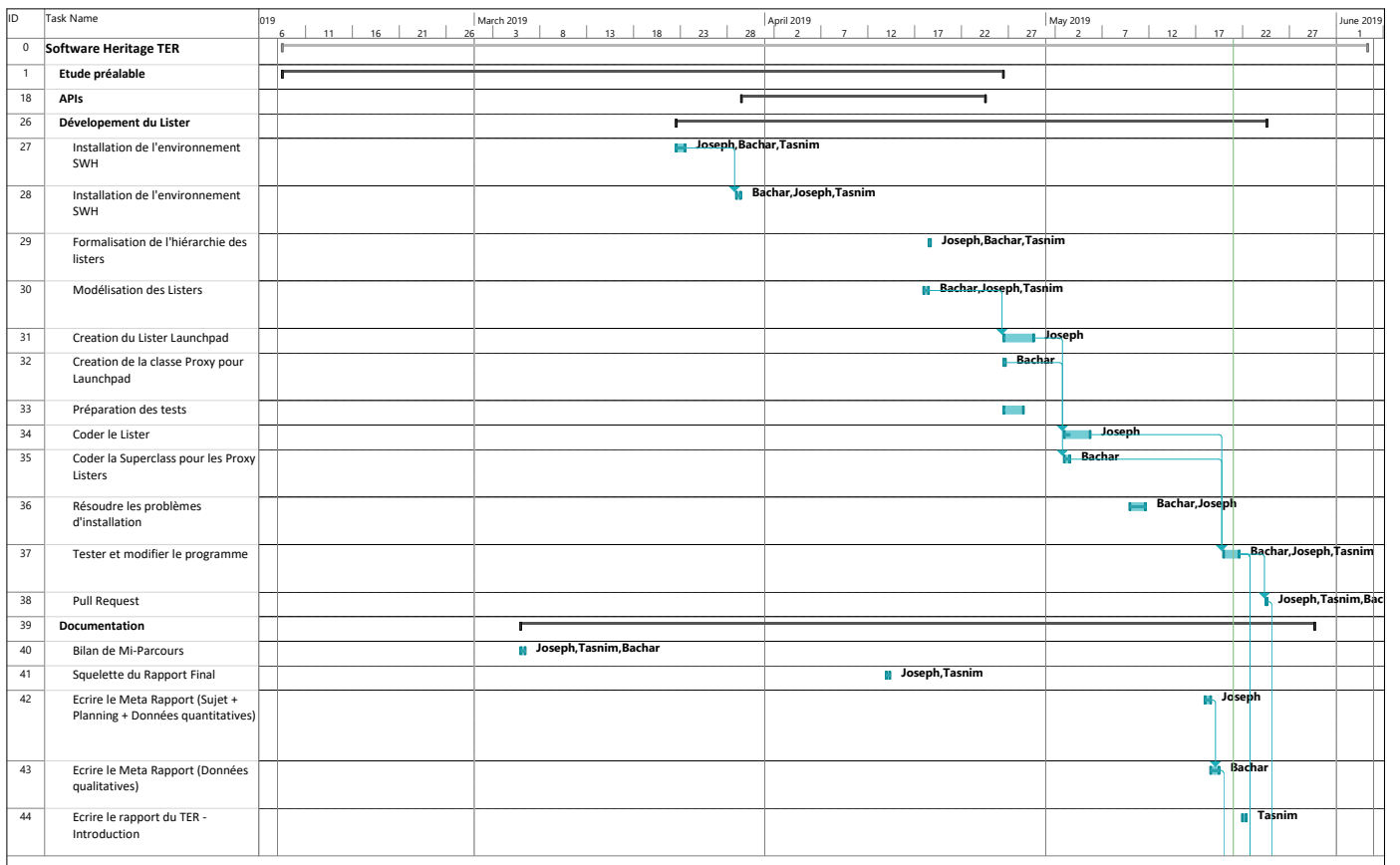
## Résultats

pull request ?

# Chapitre 7

## Conclusion

### 7.1 Planning final



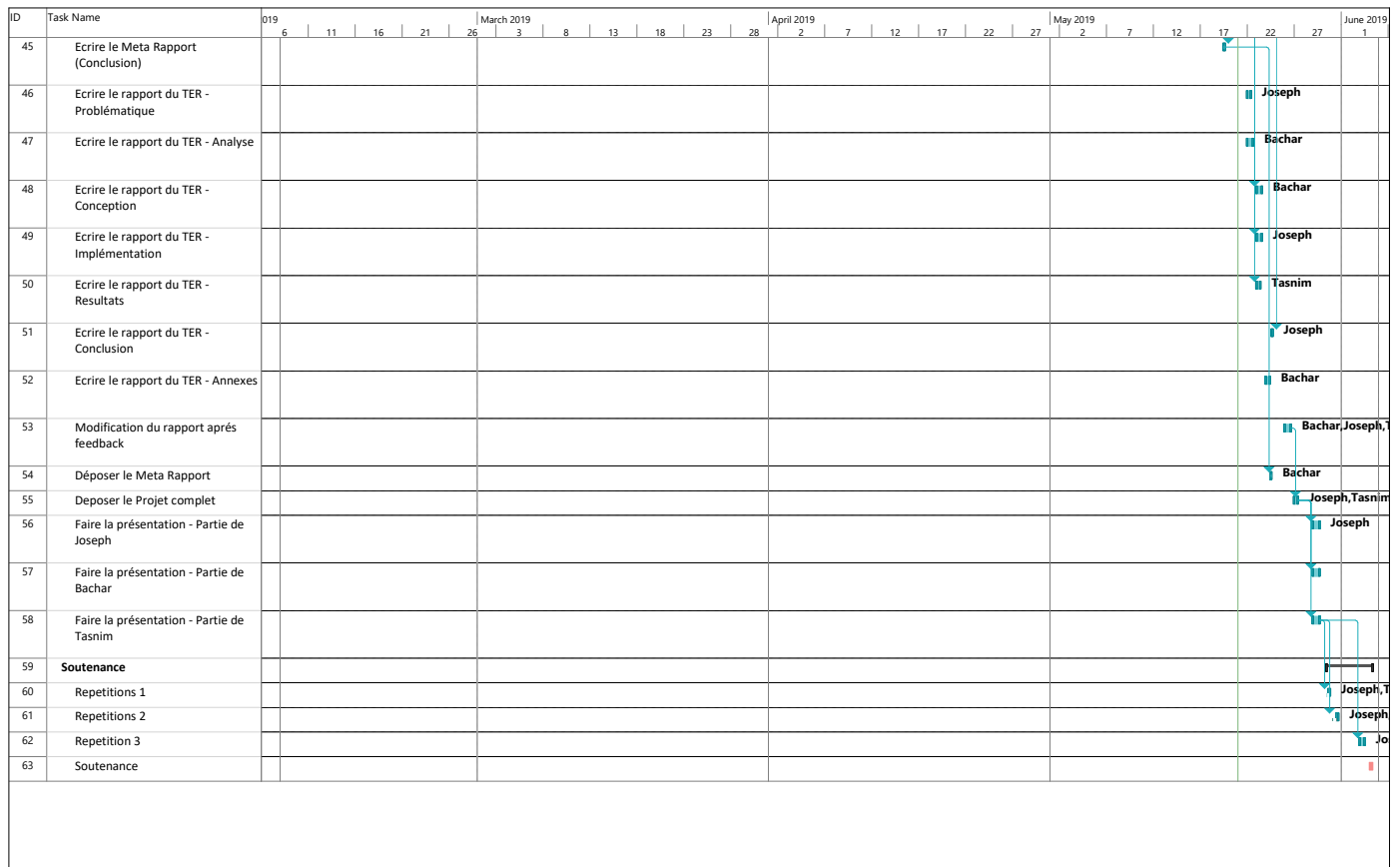


FIGURE 7.1 – Planning final

## 7.2 Difficultés rencontrées

Au cours de se projet, nous avons rencontrer des difficultés auxquelles nous ne nous attendions pas.

**7.2.1 Les APIs**

**7.2.2 Les tests**

**7.3 Perspectives**

**7.4 Bilan et apports du TER**

annexes résumés code



# Bibliographie

1)Software Heritage : Why and How to Preserve Software Source Code  
- Roberto Di Cosmo, Stefano Zacchioli 2) Identifiers for Digital Objects :  
the Case of Software Source Code Preservation - Roberto Di Cosmo, Morane  
Gruenpeter, Stefano Zacchirol