



**Sorbonne Université**

Faculté des Sciences et Ingénierie

Parcours Science et Technologie du Logiciel  
(STL)

# **Rapport de projet : Benchmarking de solutions optimistes pour génération de données test à partir de JSON-Schema**

## **Auteurs**

Abdelkader Boumessaoud  
Zaky Abdellaoui

## **Encadrant**

Mohamed-Amine Baazizi  
Lyes Attouche

# Table des matières

<b>1</b>	<b>Rapport mi-parcours</b>	<b>2</b>
1.1	contexte du projet . . . . .	2
1.2	cahier des charges . . . . .	2
1.3	tâches déjà réalisées . . . . .	2
1.4	prévisionnel des tâches restantes . . . . .	2
<b>2</b>	<b>Questions</b>	<b>3</b>
<b>3</b>	<b>Introduction</b>	<b>3</b>
3.1	Contexte et Motivation . . . . .	3
3.2	But et Objectives . . . . .	3
3.3	Étendue du projet . . . . .	3
<b>4</b>	<b>L'état de l'art</b>	<b>3</b>
4.1	Vue d'ensemble de JSON-Schema . . . . .	3
4.1.1	Exemple 1 : Générale . . . . .	4
4.1.2	Exemple 2 : Objet . . . . .	5
4.1.3	Exemple 3 : Array . . . . .	7
4.1.4	Exemple 4 : Composition . . . . .	8
4.3	Limitations des approches existantes (non-optimiste) . . . . .	10
4.3.1	limitations théoriques . . . . .	10
4.3.2	limitations pratiques . . . . .	10
<b>5</b>	<b>Méthodologie</b>	<b>10</b>
5.1	json-everything . . . . .	10
5.2	json-schema-faker . . . . .	10
5.2.1	quelques tests . . . . .	10
<b>6</b>	<b>Résultats</b>	<b>11</b>
6.1	Trouvailles . . . . .	11
6.2	Comparaison des bibliothèques . . . . .	11
6.3	Limites des approches optimistes . . . . .	11
6.3.1	limitations théoriques . . . . .	11
6.3.2	limitations pratiques . . . . .	11
<b>7</b>	<b>Conclusions</b>	<b>11</b>
7.1	Résumé des résultats . . . . .	11
7.2	Implications . . . . .	11
7.3	Travail futur . . . . .	11

# 1 Rapport mi-parcours

## strucutre

1. la description du contexte du projet
2. la présentation du cahier des charges (tâches à réaliser)
3. l'explication des tâches déjà réalisées
4. un prévisionnel des tâches restantes, avec retro-planning associé

### 1.1 contexte du projet

- objectif du projet : caracteriser les limitations des générateurs open-source en terme de classes de schémas traitées correctement et de classe de schémas problématique
- expliquer contexte : JSON Schema et
- probleme de génération d'exemples : pourquoi c'est compliqué
- citer approaches existants et outils en ligne
- dire pouquoi se limiter au Faker et JE

### 1.2 cahier des charges

### 1.3 tâches déjà réalisées

- étude et compréhension : syntaxe et sémantique du langage
- prise en main d'un validateur : experimentation sur Test Suite officiel [?], familiarisation avec sortie standard
- prise en main de générateurs open-source [?] et [?] : expérimentations sur collection X (xx schémas) et Y (yy schémas)
- implantation d'une chaine de traitement : générateur + validateur avec sortie détaillée
- analyse du code guidée par les exemples : une fois une class de schéma provoquant des erreurs identifiée, synthétiser des schémas ayant les mêmes caracteristiqus (opéretueurs, profondeur, etc) et tester sur ces schémas pour confirmer l'hypothese sur une eventuelle limitation par rapport à une class de schémas
- expérimentation à petite échelle : prendre WP ou SP, classer par famille d'approches, extraire les erreurs au format standard, examiner le premier niveau seulement puis tenter de généraliser. Autre possibilité, examiner pour chauqe output les erreurs de façon complete puis corrélér au code source pour confirmer certaines hypotheses sur la limitation de l'outil pour certains opérateurs ou agencement d'opérateurs

### 1.4 prévisionnel des tâches restantes

- experimentations à large échelle : considérer les collections de schémas de plus grande taille, confirmer/infirmir les hypotheses dérivées sur les collections de plus petite taille

## 2 Questions

- demander si on peut checker le travail fait sur la 3eme lib - discuter de notre plan d'organisation entre le faker et le everything

## 3 Introduction

### 3.1 Contexte et Motivation

JSON-Schema continue d'être un outil précieux pour définir et valider la structure des données JSON dans une variété de contextes. Son expressivité et son format standardisé en font un choix idéal pour décrire les API, vérifier la structure des données dans les pipelines d'apprentissage et spécifier les exigences du développement logiciel. L'utilisation de JSON-Schema dans ces domaines a augmenté ces dernières années, ce qui rend important de disposer d'outils efficaces pour générer des instances de test à partir de ces schémas.

De plus, l'essor de l'architecture des microservices et le besoin d'interopérabilité des données ont augmenté la demande de JSON-Schema, car il permet la définition de structures de données partagées entre plusieurs services. Son utilisation a également été étendue à d'autres applications telles que la validation de données dans des bases de données et des systèmes de stockage de données.

Dans l'ensemble, l'utilité de JSON-Schema découle de sa capacité à fournir un moyen standardisé et expressif de définir et de valider la structure des données JSON, ce qui en fait un outil indispensable dans de nombreux domaines du développement logiciel et de la gestion des données.

Mais malgré cela, ses limites sont à dénombrer, que ce soit du point de vue de complexité, manque de normalisation, types de données limités, performance, Absence de transformation des données, rapport d'erreur limité, Spécification incomplète.

### 3.2 But et Objectives

L'objectif du projet est d'étudier les limitations (théoriques et pratiques) de différentes implémentations visant à générer des instances de test à partir de schémas JSON, en se concentrant principalement sur trois bibliothèques open-source. L'étude utilisera à la fois la retro-ingénierie et l'analyse expérimentale pour évaluer l'efficacité et les limites de ces approches, et pourra s'étendre à d'autres bibliothèques open source jugées applicables.

### 3.3 Étendue du projet

L'étude de ses limites permettra de les surmonter, ce qui mènera à l'élargissement de son domaine applicatif, Cela aiderait à garantir la qualité des instances de test générées à partir des schémas d'entrée, conduisant à des résultats plus fiables et plus précis. La possibilité de générer automatiquement des instances de test à partir de schémas améliorera l'efficacité des processus de développement logiciel, facilitant ainsi le test et la validation de nouvelles applications et de nouveaux systèmes.

## 4 L'état de l'art

### 4.1 Vue d'ensemble de JSON-Schema

JSON-Schema est un format de validation pour les données JSON. Il définit la structure attendue d'un document JSON en spécifiant les types de données pour chaque propriété, ainsi que les règles de validation telles que les valeurs minimales et maximales autorisées et les attributs obligatoires. Les schémas JSON peuvent être utilisés pour valider les données à la volée, ce qui permet de garantir la qualité et la cohérence des données et peuvent également être partagés entre différentes applications pour faciliter la communication de données structurées.

#### 4.1.1 Exemple 1 : Générale

Voici un exemple pour illustrer la structure de base d'un JSON-schema :

```
1 {
2   "$schema": "http://json-schema.org/draft-07/schema#",
3   "title": "Product",
4   "description": "A product from Acme's catalog",
5   "type": "object",
6   "properties": {
7     "productId": {
8       "description": "The unique identifier for a product",
9       "type": "integer",
10      "minimum": 1
11    },
12    "productName": {
13      "description": "Name of the product",
14      "type": "string",
15      "minLength": 1
16    },
17    "price": {
18      "description": "The price of the product",
19      "type": "number",
20      "exclusiveMinimum": 0
21    },
22    "tags": {
23      "description": "Tags associated with the product",
24      "type": "array",
25      "items": {
26        "type": "string"
27      },
28      "minItems": 1,
29      "uniqueItems": true
30    },
31    "dimensions": {
32      "description": "The dimensions of the product",
33      "type": "object",
34      "properties": {
35        "length": {
36          "type": "number"
37        },
38        "width": {
39          "type": "number"
40        },
41        "height": {
42          "type": "number"
43        }
44      },
45      "required": ["length", "width", "height"]
46    }
47  },
48  "required": ["productId", "productName", "price"]
49 }
```

Ce schéma décrit un objet JSON "Product" qui représente un produit dans le catalogue d'Acme. Il a les propriétés suivantes :

- "productId" (obligatoire) est un entier unique qui représente l'identifiant du produit, avec une valeur minimale de 1.
- "productName" (obligatoire) est une chaîne de caractères qui représente le nom du produit, avec une longueur minimale de 1 caractère.
- "price" (obligatoire) est un nombre qui représente le prix du produit, avec une valeur minimale ex-

clusive de 0.

- "tags" est un tableau de chaînes de caractères qui représentent les tags associés au produit, avec au moins 1 élément et des éléments uniques.
- "dimensions" est un objet qui représente les dimensions du produit, avec les propriétés "length", "width" et "height", toutes obligatoires et de type nombre.

De plus, les propriétés "productId", "productName" et "price" sont requises pour que l'objet soit considéré comme valide.

Un document JSON validé par ce schéma serait :

```
1 {
2   "productId": 42,
3   "productName": "Super Product",
4   "price": 19.99,
5   "tags": ["awesome", "must-have"],
6   "dimensions":
7     {
8       "length": 10,
9       "width": 5,
10      "height": 2
11    }
12 }
```

Un document JSON non validé ce schéma serait :

```
1 {
2   "productId": 0,
3   "productName": "",
4   "price": -1,
5   "tags": [],
6   "dimensions": {}
7 }
```

Ce document ne serait pas validé car :

- "productId" doit être un entier avec une valeur minimale de 1, mais il est défini sur 0.
- "productName" doit être une chaîne non vide, mais c'est une chaîne vide.
- "price" doit être un nombre non négatif, mais il est défini sur -1.
- "tags" doit avoir au moins 1 élément et tous les éléments doivent être uniques, mais c'est un tableau vide.
- "dimensions" doit avoir les propriétés "longueur", "largeur" et "hauteur", mais c'est un objet vide.

Le schéma requiert que toutes les propriétés "productId", "productName" et "price" soient présentes, mais ce document ne les a pas toutes.

#### 4.1.2 Exemple 2 : Objet

Ce schéma JSON décrit les propriétés et les règles pour un objet de produit :

```
1 {
2   "$schema": "http://json-schema.org/draft-07/schema#",
3   "title": "Objet de produit",
4   "description": "Un schéma pour décrire un objet de produit",
5   "type": "object",
6   "properties": {
7     "id": {
8       "type": "integer",
9       "description": "Identifiant unique du produit",
10      "minimum": 1
11    },
12  }
```

```

12     "nom": {
13         "type": "string",
14         "description": "Nom du produit"
15     },
16     "prix": {
17         "type": "number",
18         "description": "Prix du produit",
19         "minimum": 0
20     },
21     "disponible": {
22         "type": "boolean",
23         "description": "Indique si le produit est disponible"
24     }
25 },
26 "patternProperties": {
27     "^categories.*$": {
28         "type": "string",
29         "description": "Catégories du produit"
30     }
31 },
32 "additionalProperties": true,
33 "required": [
34     "id",
35     "nom",
36     "prix"
37 ],
38 "propertyNames": {
39     "pattern": "^[a-zA-Z0-9]+$",
40     "description": "Les noms de propriété ne doivent contenir que des
    caractères alphanumériques"
41 }
42 }

```

- patternProperties définit un modèle pour les propriétés qui doivent correspondre au modèle categories.\*\$. - additionalProperties est défini sur false, ce qui signifie qu'aucune propriété supplémentaire ne peut être ajoutée à l'objet de produit. - unevaluatedProperties définit un autre schéma pour les propriétés qui ne sont pas incluses dans les propriétés requises. - required définit les propriétés requises pour un objet de produit. - propertyNames définit un modèle pour les noms de propriété.

Un document JSON validé par ce schéma serait :

```

1  {
2  {
3      "id": 1,
4      "nom": "Ordinateur portable",
5      "prix": 500,
6      "disponible": true,
7      "categories": "informatique, ordinateur",
8      "unevaluatedProperties": {
9          "note": "Bonne performance pour le prix"
10     }
11 }

```

Ce document

Un document JSON non validé par ce schéma serait :

```

1  {
2      "id": -1,
3      "nom": "Ordinateur portable",
4      "prix": -500,
5      "disponible": true,

```

```

6   "categories": ["informatique", "ordinateur"],
7   "unevaluatedProperties": {
8     "note": 123
9   },
10  "taille": 200
11 }

```

Ce document n'est pas valide pour les raisons suivantes :

- La propriété id a une valeur négative, ce qui est en dehors de la plage minimale définie dans le schéma.
- La propriété prix a une valeur négative, ce qui est en dehors de la plage minimale définie dans le schéma.
- La propriété unevaluatedProperties.note n'est pas une chaîne de caractères, ce qui est en dehors des types définis dans le schéma.
- La propriété taille est mal orthographiée dans le document JSON, ce qui signifie qu'elle n'est pas définie dans le schéma et que additionalProperties est défini sur false.

#### 4.1.3 Exemple 3 : Array

Voici un schéma JSON qui définit un tableau représentant des produits :

```

1  {
2    "$schema": "http://json-schema.org/draft-07/schema#",
3    "type": "array",
4    "items": [
5      {
6        "type": "string",
7        "minLength": 3,
8        "maxLength": 10
9      },
10     {
11       "type": "integer",
12       "minimum": 0
13     }
14   ],
15   "tuple": [
16     {
17       "type": "string"
18     },
19     {
20       "type": "integer",
21       "minimum": 10
22     }
23   ],
24   "contains": {
25     "type": "string",
26     "enum": ["apple", "banana", "cherry"]
27   },
28   "minItems": 2,
29   "uniqueItems": true
30 }

```

Ce schéma définit un tableau qui doit avoir au moins 2 éléments et dont les éléments sont uniques. Les éléments du tableau doivent être soit une chaîne avec une longueur minimale de 3 et une longueur maximale de 10, soit un nombre entier supérieur ou égal à 0. Le tableau doit également contenir une séquence de deux éléments, une chaîne et un nombre entier supérieur ou égal à 10. Enfin, le tableau doit contenir au moins un élément parmi une liste d'options prédéfinies (pomme, banane, cerise).

- Items : Il s'agit d'une propriété qui définit les critères de validation pour les éléments d'un tableau. Il peut s'agir d'un schéma JSON complet ou simplement d'un type de données (comme "integer" ou "string").



- Tuple validation : La propriété "items" peut également être définie comme une liste de schémas JSON, permettant une validation tuple pour le tableau, ce qui signifie que chaque élément du tableau doit respecter le schéma défini à la même position dans la liste.

- Unevaluated items : La propriété "additionalItems" peut être utilisée pour définir comment les éléments supplémentaires dans le tableau doivent être traités si le nombre d'éléments dépasse la longueur de la liste de schémas définis dans "items". Si "additionalItems" est défini comme "false", les éléments supplémentaires seront considérés comme non valides. Si "additionalItems" est défini comme un schéma JSON, les éléments supplémentaires seront validés contre ce schéma.

- Contains : La propriété "contains" peut être utilisée pour définir un schéma JSON qui doit être inclus dans au moins un élément du tableau.

- Length : Les propriétés "minItems" et "maxItems" peuvent être utilisées pour définir une plage minimale et maximale de nombres d'éléments pour le tableau.

- Uniqueness : La propriété "uniqueItems" peut être définie comme "true" pour exiger que tous les éléments du tableau soient uniques.

Un document JSON validé par ce schéma serait :

```
1  [
2    "apple", 10,
3    "banana", 20
4  ]
```

Un document JSON non validé par ce schéma serait :

```
1  [
2    "est",
3    5979550
4  ]
```

Ce document (généré par le faker) ne serait pas valide car : ne respecte pas le contains

#### 4.1.4 Exemple 4 : Composition

Voici un schéma JSON qui utilise la composition de schémas :

```
1  {
2    "$schema": "http://json-schema.org/draft-07/schema#",
3    "title": "Personne",
4    "type": "object",
5    "properties": {
6      "nom": {
7        "type": "string",
8        "minLength": 3
9      },
10     "age": {
11       "type": "integer",
12       "minimum": 0
13     },
14     "sexe": {
15       "type": "string",
16       "enum": ["homme", "femme"]
17     }
18   },
19   "required": ["nom", "age", "sexe"],
20   "allOf": [
21     {
22       "not": {
23         "properties": {
24           "age": {
```

```

25         "maximum": 150
26     }
27 }
28 }
29 },
30 {
31     "anyOf": [
32     {
33         "properties": {
34             "sexe": {
35                 "const": "homme"
36             }
37         }
38     },
39     {
40         "properties": {
41             "age": {
42                 "minimum": 18
43             }
44         }
45     }
46 ]
47 },
48 {
49     "oneOf": [
50     {
51         "properties": {
52             "nom": {
53                 "pattern": "^[A-Z] [a-z]+$"
54             }
55         }
56     },
57     {
58         "properties": {
59             "nom": {
60                 "pattern": "^[A-Z] [a-z]+ [A-Z] [a-z]+$"
61             }
62         }
63     }
64 ]
65 }
66 ]
67 }

```

Ce schéma définit un objet Personne qui a des propriétés obligatoires nom, age et sexe.

L'opérateur `allOf` spécifie que toutes les conditions doivent être respectées. La première condition utilise l'opérateur `not` pour spécifier que l'âge ne doit pas dépasser 150 ans. La deuxième condition utilise `anyOf` pour spécifier que le sexe doit être "homme" ou que l'âge doit être d'au moins 18 ans. La troisième condition utilise `oneOf` pour spécifier que le nom doit être soit en format "Prénom Nom" soit en format "Nom".

Un document JSON validé par ce schéma serait :

```

1  {
2  {
3      "nom": "John Doe",
4      "age": 170,
5      "sexe": "homme"
6  }

```

Cet objet a un nom qui est soit "John Doe" en suivant le pattern

La génération d'instances de tests à partir de JSON-Schema consiste à créer des exemples de données qui respectent la structure définie par le schéma JSON. Cela peut être utile pour les tests automatisés, la validation des entrées utilisateur, etc. Les outils de génération d'instances de tests peuvent utiliser les contraintes définies dans le schéma, telles que les types de données, les tailles minimales et maximales, les valeurs minimales et maximales, les éléments uniques, etc., pour générer des données qui respectent ces contraintes.

Les techniques existent pour générer automatiquement des instances de tests à partir de schémas JSON ne garantissent pas toujours la conformité. Une approche optimiste est de combiner des fragments de schémas pour produire une instance valide avec l'utilisation d'un validateur externe pour corriger les instances non conformes. Son efficacité est dans l'absence de la complexité difficile du problème de génération, qui a une complexité exponentielle prouvée.

### 4.3 Limitations des approches existantes (non-optimiste)

Il existe plusieurs approches pour générer des instances de test à partir de JSON-Schema :

- Génération aléatoire : en utilisant des bibliothèques telles que jsonschema2popo, on peut générer des instances aléatoires en se basant sur le JSON-Schema, en respectant les contraintes définies dans le schéma.
- Génération basée sur des exemples : en utilisant des bibliothèques telles que json-schema-faker, on peut générer des instances en se basant sur des exemples spécifiés dans le JSON-Schema.
- Génération à partir de modèles : en utilisant des outils tels que QuickType, on peut générer du code pour des instances de test en fonction du JSON-Schema, ce qui peut être utile pour les tests automatisés.

#### 4.3.1 limitations théoriques

#### 4.3.2 limitations pratiques

## 5 Méthodologie

### 5.1 json-everything

### 5.2 json-schema-faker

#### 5.2.1 quelques tests

les instances générés par json-schema-faker (sur interface) seront testés sur jsonschemavalidator en prenant l'exemple 1 : objet

## 6 Résultats

### 6.1 Trouvailles

### 6.2 Comparaison des bibliothèques

### 6.3 Limites des approches optimistes

#### 6.3.1 limitations théoriques

#### 6.3.2 limitations pratiques

## 7 Conclusions

### 7.1 Résumé des résultats

### 7.2 Implications

### 7.3 Travail futur

## Références

- [1] <https://json-schema.org>
- [2] Jsongenerator. <https://github.com/jimblackler/jsongenerator>
- [3] Json schema faker. <https://github.com/json-schema-faker/json-schema-fakeryes>
- [4] JSON everything. <https://github.com/gregsdennis/json-everything>
- [5] Lyes Attouche, Mohamed Amine Baazizi, Dario Colazzo, Francesco Falleni, Giorgio Ghelli, Cristiano Landi, Carlo Sartiani, Stefanie Scherzinger : A Tool for JSON Schema Witness Generation. EDBT 2021 : 694-697
- [6] Mohamed Amine Baazizi, Dario Colazzo, Giorgio Ghelli, Carlo Sartiani, Stefanie Scherzinger : Not Elimination and Witness Generation for JSON Schema. CoRR abs/2104.14828 (2021)
- [7] Pierre Bourhis, Juan L. Reutter, Fernando Suárez, Domagoj Vrgoc