

## コード自動生成（CodeGeeX 他）

### 目次

はじめに .....	1
1. コードの自動生成が生まれた経緯.....	1
2. AI を使うリスク.....	2
2.1. 著作権等 .....	2
2.2. 漏洩.....	2
2.3. 生成されるコードの精度.....	2
3. コード生成の範囲やタイプ.....	3
4. コード生成ツールの試用 .....	3
4.1. CodeGeeX について（注意点） .....	3
4.2. CodeGeeX のインストール.....	4
4.3. CodeGeeX の使用 .....	5
4.3.1. 英文のコメントからコード生成.....	5
4.3.2. 日本語のコメントからコード生成.....	7
4.3.3. メソッドコメントの書き方による変化.....	9
4.3.4. 処理内容からコメントを生成 .....	11
5. コード自動生成のその他のツール.....	12
5.1. ブラウザで試せる Codeium.....	12
5.2. オープンソースの言語モデルを利用する Hugging Face.....	13
5.3. スタックオーバーフロー検索（CROKAGE、Que2Code 他） .....	13
6. 現状の評価 .....	13
7. 今後 .....	14

## コード自動生成 (CodeGeeX 他)

はじめに

プログラミングに付いて回る決まりきったコードを減らす取り組みは汎用機 (メインフレーム) のマクロアセンブラの時代からありました。よく使われるものの一つに画面生成マクロがあり、項目毎に表示位置や入出力、文字種類、桁数等の固有情報を指定するだけで業務画面が完成します。

現代ではフレームワークや共用ライブラリを使うことが定番になり、制御 (IO や判定、繰り返し) のコードを削減する Java の Stream や C# の Linq、Python の内包表記、ラムダ式、…等の簡明な書き方を言語仕様に取り込み続けています。また、SQL クエリを使えば論理的な条件を記述するだけで制御コードを書かずにデータ間の関連付けや加工ができます。

次に向かおうとしているのは Auto Coding や Code Generator と呼ぶコードの自動生成です。これはインターネット上の大量の情報を学習データとし、大量の CPU/GPU とメモリ、ストレージを使って処理・作成した事前学習済言語モデルからコードを生成しようという新しいアプローチで、それまでのテンプレートに変動パラメータを埋め込むソース生成とは全く異なるものです。

### 1. コードの自動生成が生まれた経緯

2018 年に Google 社が自然言語処理<sup>1</sup>の論文 BERT(Bidirectional Encoder Representations from Transformers : Transformer を活用した双方向的エンコード表現)とその実装を発表したのを端緒に、2019 年に OpenAI 社が Transformer をベースにインターネット上の大量の情報で事前学習した言語モデル GPT-2(Generative Pre-trained Transformer-2)を公開して自然な会話ができる AI<sup>2</sup>として話題になりました。2020 年に次版の GPT-3 の独占的ライセンス契約を Microsoft 社と結び Microsoft 社はこれを使って無料の開発ツール CodeGPT (開発環境 VSCode の機能拡張) の提供を始めました。BERT も GPT もテキストを変換 (Transform : 翻訳、要約、他言語生成、...) のための言語モデルとしてつくられており、開発ツールとして使うためには変換前後のテキストを開発環境で取り扱う機能が必要で CodeGPT がこれを行います。2023 年 4 月現在 CodeGPT は以下の機能を持っています。

- ・ GPT-4 へのアクセス
- ・ ChatGPT とのチャット
- ・ コメントからコードを生成
- ・ ソースコードのリファクタリング
- ・ ソースコードからコメントを生成
- ・ その他

CodeGPT は無償で使えますが GPT-3 へのアクセスには OpenAI 社への登録が必要なため、色々な組織が Transformer ベースの事前学習済言語モデル (BERT、GPT 以外にも T5、DeBERTa 等々) と言語モデルを使うためのツールを競って開発しており Auto-Coding や Code-Generator 等のタグを付けてオープンソースで公開しています。因みに GPT-2 もオープンソース (修正 MIT ライセンス) で公開<sup>3</sup>されています。

---

<sup>1</sup> 言語処理の方法は「bert gpt attention」等の文言で検索すると解説サイトが沢山でてきます

<sup>2</sup> GPT-2,GPT-3,GPT-4 は大量の学習データを使い言語データを整理した事前学習済言語モデルです  
一般に AI と呼ばれている ChatGPT は GPT-3.5 に会話インタフェースを付加したものです

<sup>3</sup> GitHub/openai/gpt-2 <https://github.com/openai/gpt-2>

## コード自動生成 (CodeGeeX 他)

### 2. AI を使うリスク

AI 全般に関して SF 的な漠然とした不安を語る論調がありますが、コードの自動生成に関しては黎明期特有のもう少し具体的なリスクがあります。

#### 2.1. 著作権等

コード生成は情報源としてコード生成用の言語モデルを使います。生成用の言語モデルは GitHub のリポジトリ等に登録されたソースコードから事前学習したものが多くありますが、各ソースには著作権があり、コードによっては特許申請・登録されたものが混じっているかもしれません。

日本の著作権法では学習データとして使う場合は著作権者の承諾は不要のよう<sup>4</sup>ですが、AI 利用の切っ掛けになった GPT-2 からは学習データの一部を抽出することができるという報告<sup>5</sup>があり、実際にこれを行うツールが GitHub で公開されています。また、自動生成ツールが提案したコードが事前学習データに使った商用禁止ライセンスのソースと完全一致した場合、偶然の一致と主張できるか否か (小さく分割されたトークンから問合せ内容に基づいて新たに組み立てているので問題なさそうな気はしますが...)、海外で運用・利用されているサイトのソースで作った言語モデルから自動生成したソースの著作権の帰属先がどう判定されるのか、現状では分かりません。

大半の言語モデルは事前学習に使用するソースのライセンスに言及していませんが、モデルによってはライセンスフリーのソースだけを使っていると宣言しているものもあります。

#### 2.2. 漏洩

開発環境に組み込んでコードの自動補完を受けるためにエディターで編集中のコードの全体か一部を言語モデルに送りますが、送られたコードは言語処理をするために事前学習データと同一の機能を使ってエンコードされ言語モデルに取り込まれる (学習データとして使われる) 可能性があります。ツールによっては取込みを停止することもできますが、IP アドレスを含むアクセス記録は残ります。

また、CodeGPT やその他の代替ツールはリファクタリングの機能を持っています。この機能を使うためには自分で開発しているソースを言語モデル側に送る必要がありますが、転送経路や送った先で保護される保証はありませんし、言語モデルに取り込まれる可能性は自動補完と同様です。

#### 2.3. 生成されるコードの精度

コード生成で使われる言語モデルは各種のプログラミング言語用に調整 (ファインチューニング) されていると推測しますが、特定のバージョンで追加されたり廃止されたクラスや記法をどのように適用するのか明示されたものは見当たりません。また、言語モデルは既に誰かが記述した蓄積を基にしたものなので、調整次第ですが、より新しい機能 (クラス、メソッド) やアルゴリズムではなく最大公約数的なソースの提案になる可能性が強くなります。更に「ChatGPT に生成させたコードは生成する回答が間違っている可能性が高い」ので使用を当面禁止としたコミュニティもあります。

---

<sup>4</sup> 総務省「AI を用いたクラウドサービスに関するガイドブック」4.1.2 留意すべき事項 3) データ収集と著作権 [https://www.soumu.go.jp/menu\\_news/s-news/01ryutsu06\\_02000305.html](https://www.soumu.go.jp/menu_news/s-news/01ryutsu06_02000305.html)

<sup>5</sup> Extracting Training Data from Large Language Models <https://arxiv.org/abs/2012.07805>

## コード自動生成（CodeGeeX 他）

### 3. コード生成の範囲やタイプ

コード生成／コード補正を行う各種のツールは、それぞれ以下のレベルを目標にしています。

- ① クラス全体を生成する
- ② Spring や FLASK、React 等のフレームワークに組み込み可能なコードを自動生成する
- ③ メソッドや処理ブロックを生成する
- ④ 入力中の文や式の補完し、次の文や式を提案する

①～③のようにある程度高度な機能を実現するにはツールに仕様を伝える方法が難しくなります。

③は連続する複数行のコメントで実現している例もありますが、決まりきった書き方を覚える必要があるのであれば直接 コードを覚えた方が効率がよいかもしれません。

①、②が満足できるレベルで実現できるようになればプログラミング言語を習得する必要がなくなり是非実現して貰いたいですが、現状は④のレベルが現実的で、しかも言語モデルを処理するために開発環境として必要なマシン性能もかなり高いものになります。

### 4. コード生成ツールの試用

2023 年 4 月時点で実際に動作するコード自動生成のツール（言語モデル×開発環境のツール）には、有償（ツール or/and 言語モデル）、今だけ無償で利用者登録要、等のものがありますが、本資料では無償で利用者登録不要の CodeGeeX を使ってコード生成のイメージを確認します。

#### 4.1. CodeGeeX について（注意点）

CodeGeeX は GitHub に登録されているオープンソースで、以下のように紹介されています<sup>6</sup>。

「CodeGeeX は、20 を超えるプログラミング言語の大規模なコード コーパスで事前トレーニングされた、130 億のパラメーターを持つ大規模な多言語コード生成モデルです。2022 年 6 月 22 日の時点で、CodeGeeX は 1,536 個の Ascend 910 AI プロセッサのクラスターで 8,500 億を超えるトークンでトレーニングされています。」

そして、「CodeGeeX にはいくつかのユニークな機能があり」として以下の記載があります。

- 多言語コード生成: CodeGeeX は、Python、C++、Java、JavaScript、Go など、<以下略>
- カスタマイズ可能なプログラミング アシスタント: CodeGeeX は、VS Code 拡張機能マーケットプレイスで無料で入手できます。コード補完、説明、要約などをサポートし、<以下略>
- オープンソースとクロスプラットフォーム: すべてのコードとモデルの重みは、研究目的で公開し<以下略>

#### <2023/4/1 現在 CodeGeeX : ライセンス条項>

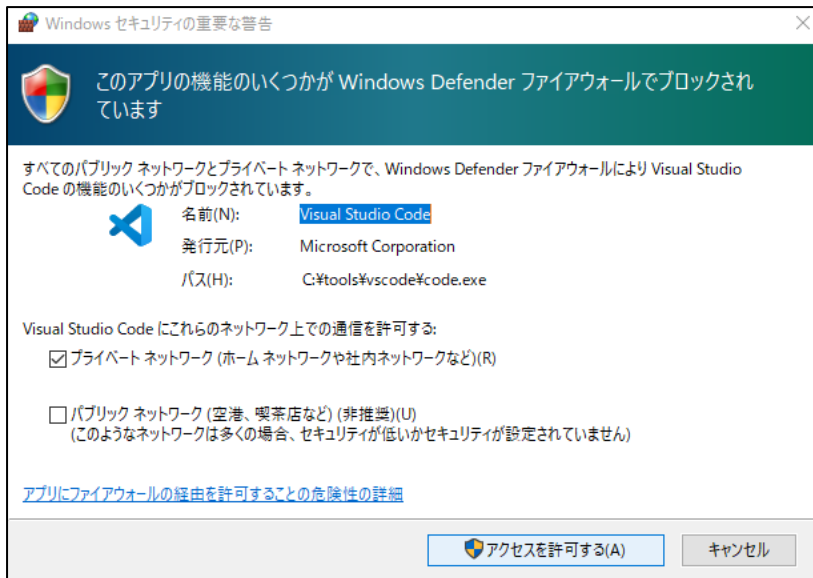
●コードは Apache-2.0 ライセンス適用で商業使用、改造、配布が許可されますが、データモデルに関しては以下の記載があり、**商用、軍用での使用と中国の弱体化に資する目的で使うことを禁じています。** [https://github.com/THUDM/CodeGeeX/blob/main/MODEL\\_LICENSE](https://github.com/THUDM/CodeGeeX/blob/main/MODEL_LICENSE)

---

<sup>6</sup> THUDM/CodeGeeX <https://github.com/THUDM/CodeGeeX>



## コード自動生成 (CodeGeeX 他)



### 4.3. CodeGeeX の使用

インストールするとデフォルトで CodeGeeX のコード補完が有効になっています。

以下、言語を Java にしてコードを入力していきます (アクセス修飾子等の省略可能な部分は省略)。

#### 4.3.1. 英文のコメントからコード生成

① メソッドの宣言を書いた後にコメント (3 行目) を書くとコードの提案が始まります。

```

1  class Test{
2      ... int selElementsAndSum(Map<String, Integer> map, String[] keys){
3          ... // Select elements of Map that match in keys and return sum of values.
4
5          int sum = 0; // Initialize sum to zero.
6      }

```

TAB キーで CodeGeeX の提案したコードが確定 (Tab 以外でクリア) します。1 行ずつ進めると...

```

1  class Test{
2      ... int selElementsAndSum(Map<String, Integer> map, String[] keys){
3          ... // Select elements of Map that match in keys and return sum of values.
4
5          ... int sum = 0; // Initialize sum to zero.
6          ... for(String key: keys){ // Iterate over keys.
7              ... if(map.containsKey(key)){ // Check if key is in map.
8                  ... sum += map.get(key); // If key is in map, add its value to sum.
9              ... }
10         ... }
11         ... return sum; // Return sum.
12     ... } // Close method.
13

```

上記のようなコードが完成します。

## コード自動生成 (CodeGeeX 他)

② クラス内に行コメントを記述するとコメントの続きやメソッド全体が提案されたりもします。

```

1 class Test2{
2     ...//Select elements of Map that match in keys and return sum of values.
        public static int sumValues(final Map<String, Integer> map, final Set<String> keys) {
            int sum = 0;
            for (final String key : keys) {
                final Integer value = map.get(key);
                if (value!= null) {
                    sum += value;
                }
            }
            return sum;
        }
    }
3     ...

```

提案を受け入れる場合は、行毎に TAB キーで確定していきます。

③ メソッドコメントを書き始めると、コメントの補完とメソッド宣言以降も提案してきます。  
@return の記述では、入力パラメータの設定値による分岐も提案してきます。

```

1 class Test2{
2     /**
3     * Select elements of Map that match in keys and return sum of values.
4     * @param map indata
5     * @param keys keys to be tested for match. If null, all keys are tested.
6     * @return int sum of values of matching elements. If map is empty or keys are null, ... */

```

メソッド宣言を入力すると、処理の提案が始まります。

```

class test{ Untitled-1 } class test2{ Untitled-2 }
1 class Test2{
2     /**
3     * Select elements of Map that match in keys and return sum of values.
4     * @param map indata
5     * @param keys keys to be tested for match. If null, all keys are tested.
6     * @return sum of values of matching elements. If keys is null, returns the sum of all values. If
7     */
8     public int sumValueOfMap(Map<String, Integer> map, String[] keys){
        int sum = 0; //sum of values of matching elements. If keys is null, returns the sum of all v
        if(keys == null){ //if keys is null, test all keys. If keys is non-null, return null if no mat

```

## コード自動生成 (CodeGeeX 他)

分岐した処理ブロックの途中で提案が停止 (下部ステータスバーに「No Suggestion」) したら”}”で処理ブロックを閉じると提案が続きます。

```

1  class Test2{
2  /**
3   * Select elements of Map that match in keys and return sum of values.
4   * @param map indata
5   * @param keys keys to be tested for match. If null, all keys are tested.
6   * @return sum of values of matching elements. If keys is null, returns the sum of all values. If keys is
7   */
8  public int sumValueOfMap(Map<String, Integer> map, String[] keys){
9      int sum = 0; //sum of values of matching elements. If keys is null, returns the sum of all values.
10     if(keys == null){ //if keys is null, test all keys. If keys is non-null, return null if no matching ele
11         for(String key:map.keySet()){ //for each key in map. If key is in map, add its value to sum.
12             if(map.get(key)!=null){ //if the value in map.get(key) is not null, add it to
13                 sum += map.get(key); //sum. If map.get(key) is null, the value in map.get
14             } //((key) is not in map. So, the value in
15         }
16         return sum; //return sum. If keys is null, return sum. If keys is non-null, return null if no matc
17     }else{ //keys are not null, test only keys. If keys are in map, add its value to sum. If keys are not
18         for(String key:map.keySet()){ //for each key in map. If key is in map, add its value to sum.
19             if(map.containsKey(key)){ //if key is in map, add its value to sum. If key is in map,
20                 sum += map.get(key); //sum. If key is in map, NaN is not in sum. So, NaN is
21             }
22         }
23         return sum; //return sum. If keys are not null, return sum. If keys are in map, add its value to s
24     }
25 }
26
27 import java.util.*; //import java.util.Collections; //import java.util.Collections; //import java.util.HashM
import java.util.Map; //import java.util.HashMap; //import java.util.Map; //import java.util.Set;

```

### 4.3.2. 日本語のコメントからコード生成

日本語は単語の区切りが解り難かったり動詞の活用や助詞の使い方が柔軟な分、プログラミングには使い難いのですが、日本語で書いた仕様が通らないと使いものにならないので試してみました。

① ファイルの先頭から日本語でコメント「List の keys で指定した値と一致する Map の要素を選択し、値の合計を返す。」を書くとそれらしいライブラリの import 文とテストコードが書かれ出します。

```

1  //Listのkeysで指定した値と一致する Map の要素を選択し、値の合計を返す。
2  import java.util.HashMap;
3  import java.util.Map;
4  import java.util.Scanner;
5  import java.util.function.Function; //for Integer::valueOf (from Guava library) 描述: https://codegeex.
6  import java.util.stream.Collectors; //for toString (from Guava library) 描述: https://codegeex.cn
7  import static java.util.stream.Collectors.groupingBy; //for (key, value) -> value (from Guava library)
8  import static java.util.stream.Collectors.summingInt; //for int (from Guava library) 描述: https://code
9  import static java.util.stream.IntStream.rangeClosed; //for int (from GUA library) 描述: https://
10 public class MapTest { //the class name must be MapTest (not Map<Integer, Integer> or Map<String, Integer> or M
11     //public static void main(String[] args) { //the main method must be named main (not main$Main) 创建一个名
12     //    Map<Integer, Integer> map = new HashMap<>(); //a new map (not a new HashMap) 描述: https
13     //    map.put(1, 1); //add a new entry (not replace an existing entry) 描述: https (key = 1,
14     //    map.put(2, 2); //add a new entry (not replace an existing entry) 描述: https (key = 2,

```

※ 一見日本語が解釈できていそうな感じもしますが、英字の List や Map に反応しただけかもしれません。そして、なぜテスト用のコードにしようとツールが判断したのかも分かりません



## コード自動生成 (CodeGeeX 他)

② “List の keys で指定した値”の部分で「List の keys の値と一致する Map の要素を選択し、値の合計を返す。」と変え、メソッド宣言を手書きすると、以下のように提案してきます。

```

1 //Listのkeysの値と一致する Map の要素を選択し、値の合計を返す。
2 import java.util.AbstractMap;
3 import java.util.HashMap;
4 import java.util.Map; // for AbstractMap.SimpleEntry<K,V> and HashMap<K,V> classes. //
5 import java.util.Map.Entry; // for SimpleEntry<K,V> class. // Note: SimpleEntry<K,V>
6 import java.util.Scanner; // for the function readLine(). // Note: readLine() is not part
7 import java.util.Set; // for the function keys(). // Note: keys() is not part of the JDK.
8 import java.util.TreeSet; // for the function keys(). // Note: keys() is not part of the
9
10 int sumValueOfMap(Map<String, Integer>, List<String> keys){
11     int sum = 0; // the value of the sum. // Note: sum is not part of the function. //
12     for(Entry<String, Integer> e: keys.toMap().entrySet()){ // for each entry in keys. //
13         sum += e.getValue(); // add the value of the entry. // Note: sum is not part of
14     } // for each entry in keys. // Note: // Note: // Note: // Note: //
15     return sum; // return the sum. // Note: sum is part of the function. // Note: sum is
16
17 }

```

③ 「List の keys で指定した値と一致する Map の要素を選択し、値の合計を返す。」に変えメソッド宣言を指定すると、JDK 1.8 で追加された Stream を使ってブロック内の処理を書き始めます。

```

1 //Listのkeysで指定した値と一致する Map の要素を選択し、値の合計を返す。
2 import java.util.HashMap;
3 import java.util.Map;
4 import java.util.Scanner;
5 import java.util.function.Function; // for Integer.valueOf (from Guava library) 描述: https://codegeex.
6 import java.util.stream.Collectors; // for toString (from Guava library) 描述: https://codegeex.cn
7
8 int sumValueOfMap(Map<String, Integer> map, List<String> keys){
9     return map.entrySet(). // Map.Entry<K,V>[] 型: https://codegeex.cn 描述: https://codegeex.cn 制DB: https:
10     .stream() // Stream<Map.Entry<K,V>> 型: https://codegeex.cn 描述: https://code
11     .map(e -> e.getValue().equals(Integer.valueOf(e.getKey().toString()))? e.get
12     Value(): null) // Function<V,V> 型: https://codegeex.cn 描述: https://codegeex.
13     .filter(Objects::nonNull) // Predicate<V> 型: https://codegeex.cn 描述: https://
14     .mapToInt(Integer::intValue) // Function<V,Integer> 型: https://codegeex.cn 描述:
15     .sum(); // int 型: https://codegeex.cn 描述: https://codegeex.cn 制DB: https
16 }
17 public class SumValueOfMap { // Main class: https://codegeex.cn 描述: https://codegeex.cn 制DB:

```

※ コメントの書き方というよりはコメントに合わせてメソッド宣言を書くのが有効でした。それに、メソッドの呼び出し側の都合から考えればメソッド名や引数を自動生成させることはないでしょう

## コード自動生成 (CodeGeeX 他)

### 4.3.3. メソッドコメントの書き方による変化

Stream を使った以下のコードが生成できるか試してみました。

```
-----
BigDecimal sumValueOfArray(String[] kingaku) {
    return Arrays.stream(kingaku)
        .filter(Objects::nonNull)
        .filter(s -> s.matches("¥d+¥¥.*¥d*")) //数値と"."*1 以外を除外
        .map(BigDecimal::new) //文字列 10 進数変換
        .reduce(BigDecimal.ZERO, BigDecimal::add); //集計
}
-----
```

#### ① メソッドコメントとして実行したい内容を記述

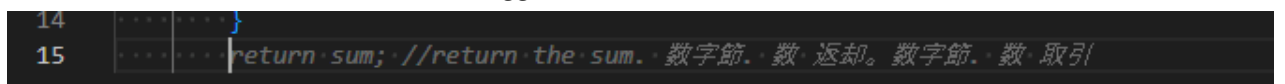


```

1  class Test4 {
2      ... /**
3      ... * 小数点付きのkingakuからStreamを作成する。
4      ... * nullと数値以外を除外する。
5      ... * Streamを10進数に変換する。
6      ... * 変換した10進数を集計する。
7      ... * @param kingaku
8      ... * @return sum
9      ... */
10     ... BigDecimal sumValueOfArray(String[] kingaku) {
11     ...     BigDecimal sum = BigDecimal.ZERO; //initial value for the sum. (初期値) 数値以外。Stream
12     ...     for (String number : kingaku) { //for each number in the array. 数字節. 数値
13     ...         sum = sum.add(new BigDecimal(number)); //add the number to the sum. 数字節. 数
14     ...     }
15     ...
16     ... }

```

ここで、下部ステータスバーに「No Suggestion」が出て止まるので、「14 行目」}」を打鍵し閉じる



```


14     ...     }
15     ...     return sum; //return the sum. 数字節. 数 返却。数字節. 数 取引

```

return が提案されて一連の処理完成します。

【結果】“null と数値以外を除外する”という部分が反映されません

#### ② コメントの句点 (「。」) を最終行だけに付けそれ以外の行から句点を削ると一気に以下の提案となりましたが、「null と数値以外を除外」はやはりコード化されていません。



```

1  class Test4 {
2      ... /**
3      ... * 小数点付きのkingakuからStreamを作成する
4      ... * nullと数値以外を除外する
5      ... * Streamを10進数に変換する
6      ... * 変換した10進数を集計する。
7      ... * @param kingaku
8      ... * @return sum
9      ... */
10     ... BigDecimal sumValueOfArray(String[] kingaku) {
11     ...     return Arrays.stream(kingaku).map(BigDecimal::new).map(BigDecimal::valueOf).sum();
12     ...
13     ... }
14 }

```

## コード自動生成 (CodeGeeX 他)

11 行目の行末にはコメントとして以下の出力がされました。

```
.sum()); //return Arrays.stream(kingaku).mapToDouble(e ->
Double.parseDouble(e)).sum(); //return Arrays.stream(kingaku).mapToLong(e ->
Long.parseLong(e)).sum(); //return Arrays.stream(kingaku).mapToInt(e ->
Integer.parseInt(e)).sum(); //return Arrays.stream(kingaku).mapToBigDecimal(e ->
new BigDecimal(e
```

③ 「null と数値以外を除外する」を「数値を選択する」に変えると、提案が変更されましたが...

```
1 class Test4 {
2     /**
3     * 小数点付きのkingakuからStreamを作成する
4     * 数値を選択する
5     * Streamを10進数に変換する
6     * 変換した10進数を集計する。
7     * @param kingaku
8     * @return sum
9     */
10    BigDecimal sumValueOfArray(String[] kingaku) {
11        return Arrays.stream(kingaku) //ArrayList<String> stream = Arrays.stream(kingaku) //stream.mapToDouble(Integer::pa
12        .....map(this::toBigDecimal) //Stream<BigDecimal> toBigDecimalStream = stream.mapToDouble(Integer:::
13        .....mapToObj(this::toString) //String to String stream = stream.mapToDouble(Integer::parseInt).mapToObj(this
14        .....map(number -> new BigDecimal(number)) //BigDecimal to BigDecimal stream = stream.mapToBigDecimal(this
15        .....reduce(BigDecimal.ZERO, BigDecimal::add) //return stream.reduce(BigDecimal.ZERO,
16        .....max(BigDecimal.ZERO, number.movePointRight(number.scale())) //return stream.mapToDouble(this::toBigDecimal);
17        .....map(number -> number.setScale(number.scale())) //return stream.mapToBigDecimal(this::toBigDecimal
18        .....map(number -> number.movePointRight(number.scale())) //return stream.mapToBigDecimal(this::toBigDec
19        .....map(number -> number.abs()) //return stream.mapToBigDecimal(this::toBigDecimal) //return stream.map
20        .....max(BigDecimal.ZERO, number.movePointRight(number.scale())) //return stream.max(this::compareTo
21        .....orElseThrow(() -> new IllegalStateException("Unknown error.")) //return stream.max(BigDecimal.ZERO).get().
22        .....toBigDecimalStream //return toBigDecimalStream.max(BigDecimal::compareTo); //return toBigDecimalStream.
23        .....reduce(BigDecimal.ZERO, BigDecimal::add); //return toBigDecimalStream.reduce(BigDecimal
24        .....add); //return toBigDecimalStream.reduce(BigDecimal.ZERO, BigDecimal::add); //return to
25        .....//BigDecimalStream.reduce(BigDecimal.ZERO, BigDecimal::add); //return toBigDecimalStream.red
```

④ コンパイルが通らなそう...

```
1 class Test4 {
2     /**
3     * 小数点付きのkingakuからStreamを作成する
4     * Streamから数字の値を選択し、10進数に変換する
5     * 変換した10進数を集計する。
6     * @param kingaku
7     * @return sum
8     */
9     BigDecimal sumValueOfArray(String[] kingaku) {
10        return new BigDecimal(kingaku.length).movePointLeft(2).map(i -> new BigDecimal(i)).reduce(BigDecimal.ZERO, BigDecimal::add);
11        .....TEN); //movePointLeft(2): a.movePointLeft(2); a.movePointLeft(2, BigDec
12        .....}
13
```

⑤ なぜか、10 進数にした後に Long⇒Integer と変換してから計算して 10 進数に戻している

```
1 class Test4 {
2     /**
3     * 小数点付きのkingakuからStreamを作成する
4     * Streamから数字の値を選択し、10進数に変換する
5     * 変換した10進数を集計して返す。
6     * @param kingaku
7     * @return sum
8     */
9     BigDecimal sumValueOfArray(String[] kingaku) {
10        return Arrays.stream(kingaku) //Stream<String> 作成(範囲限定) 或 Stream<String> 作成(終端) 以降
11        .....map(this::toBigDecimal) //Stream<BigDecimal> 作成(範囲限定) 或
12        .....mapToLong(this::toBigDecimal) //LongStream<BigDecimal> 作成(範囲限定
13        .....mapToInt(BigDecimal::intValueExact) //IntStream<BigDecimal> 作成(範囲
14        .....sum()) //BigDecimal 作成(終端)。IntStream.sum() 作成(終
15        .....toBigDecimal()); //BigDecimal 作成(終端)。toString() 作成(
16
```

## コード自動生成 (CodeGeeX 他)

### 4.3.4. 処理内容からコメントを生成

CodeGeeX のコード説明を作成する機能を実行してみました。

```

1  BigDecimal sumValueOfArray(String[] kingaku) {
2  .....
3  .....
4  ..... //数値と". "*1以外を除外
5  ..... //文字列10進数変換
6  ..... //集計
7  .....
8  }

```

CodeGeeX Tool		CodeGeeX: Explanation	
Run Code	Ctrl+Alt+N	CodeGeeX: Interactive mode	Ctrl+Enter
すべての出現箇所を変更	Ctrl+F2	CodeGeeX: Login	
リファクター...	Ctrl+Shift+R	CodeGeeX: Prompt mode(Experimental)	Alt+T
切り取り	Ctrl+X	CodeGeeX: Translation mode	Ctrl+Alt+T

使用言語の選択肢は今のところ2つだけですが...

```

//Listのkeysの値と
Choose explanation language
zh-CN
en-US

```

10 秒程度すると、2 行目が挿入されました。内容はコードと逆になってしまっていますが、一連の操作から CodeGeeX が日本語を学習したのかも (学習は拒否 disclaim を選択したはず)

```

1  BigDecimal sumValueOfArray(String[] kingaku) {
2  ..... //kingakuがnullでないものを除外して、数値と". "*1以外を除外する
3  .....
4  ..... //filter()メソッドで除外する
5  ..... //Objects::nonNullでnullを除外する
6  ..... //Objects::isNullでnullを検知する
7  ..... //map()メソッドで文字列10進数変換する
8  ..... //reduce()メソッドで集計する
9  .....
10 .....
11 ..... //数値と". "*1以外を除外
12 ..... //文字列10進数変換
13 ..... //集計
14 .....
15 }

```

再度 言語に en-US を選択して実行すると、2~7 行目が追加されました。

このコメントで基にしたコードが再現できるか否かは確認していませんが、このコメントを書くよりもコードを書いた方が手間がなさそうな粒度です。

```

1  BigDecimal sumValueOfArray(String[] kingaku) {
2  ..... //kingakuは配列である
3  ..... //filter()メソッドで除外する
4  ..... //Objects::nonNullでnullを除外する
5  ..... //Objects::isNullでnullを検知する
6  ..... //map()メソッドで文字列10進数変換する
7  ..... //reduce()メソッドで集計する
8  .....
9  .....
10 ..... //数値と". "*1以外を除外
11 ..... //文字列10進数変換
12 ..... //集計
13 .....
14 }

```

## コード自動生成（CodeGeeX 他）

### 5. コード自動生成のその他のツール

#### 5.1. ブラウザで試せる Codeium

ブラウザの機能拡張ですがインストールしなくても試すことができ、日本語も使えます。

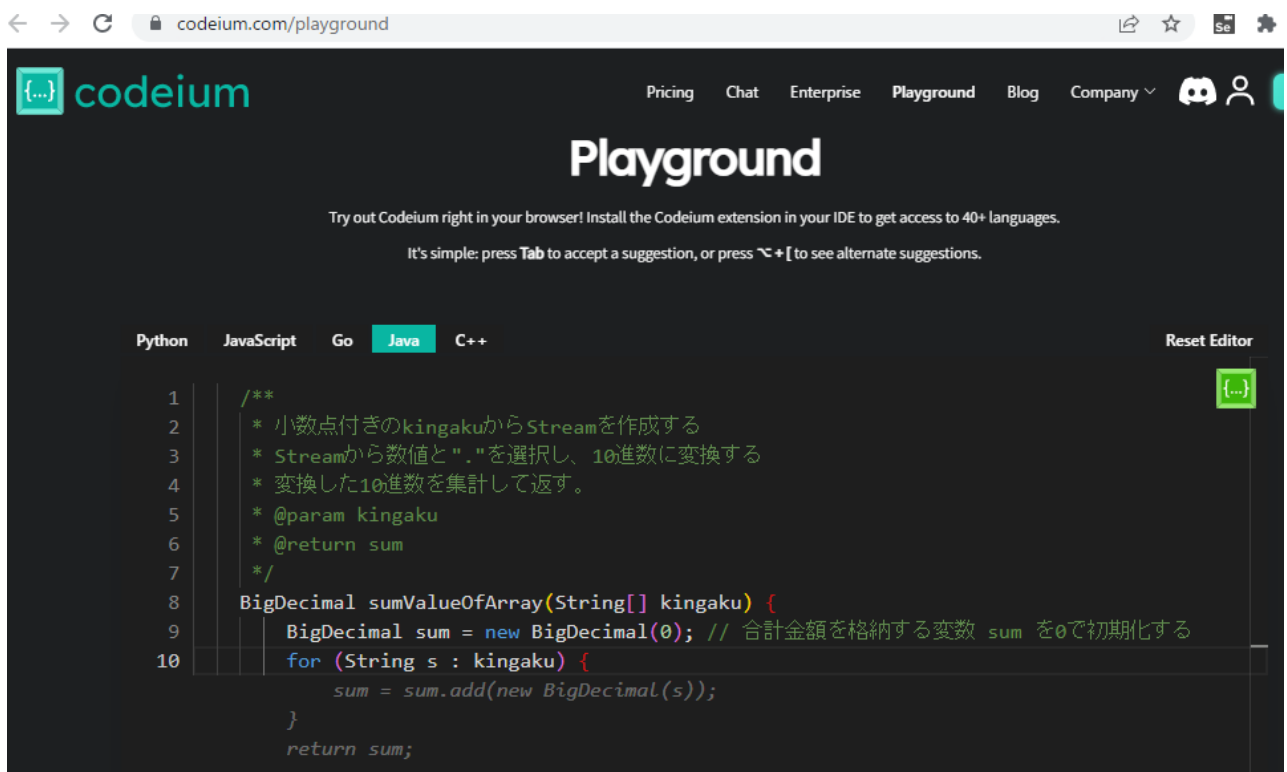
<https://codeium.com/playground>

プログラミング言語を選び、画面中央のエディター部分にコメントを入力するとコードが生成されます。Enter を押下すると 1 行ずつコードが生成され、Tab でコードが確定します。

<使用例>以下のコードを貼り付けて Enter を押下

```
/**
 * 小数点付きのkingaku から Stream を作成する
 * Stream から数値と"."を選択し、10進数に変換する
 * 変換した 10進数を集計して返す。
 * @param kingaku
 * @return sum
 */
BigDecimal sumValueOfArray(String[] kingaku) {
```

※ 9 行目以降が生成されたコードです。9 行目の行コメントも日本語で出力されました。



## コード自動生成 (CodeGeeX 他)

### 5.2. オープンソースの言語モデルを利用する Hugging Face

The AI community building the future. (未来を築く AI コミュニティ)、機械学習に関するアプリやツールを開発しているアメリカの企業です。事前学習済モデルにアクセスするための API も開発していて GitHub で公開されている BERT や GPT-2 を含むオープンソースの 18 万以上 (2023/4 時点) の事前学習済モデルがこの API を通して使えるようになっています<sup>7</sup>。

VSCoDe 用の機能拡張<sup>8</sup>は無料 (MIT ライセンス) で公開されていて、補完機能等が使えます。

### 5.3. スタックオーバーフロー検索 (CROKAGE、Que2Code 他)

スタックオーバーフロー<sup>9</sup> (Stack OverFlow) は世界中のプログラマーがコードに関するアドバイスを求めて投稿し、有識者がコード例や改善案を回答するサイトです。正式ドキュメントだけでは分からない具体的な (殆どはそのまま動作可能な) 実装例が実際の開発者から寄せられるので、現段階では「コードの自動生成」よりもこちらのコード例を参考にすることが開発の助けになるかもしれません。

このサイトに登録されているコード例 (スニペット) を探すためのツールも幾つか開発されており、CROKAGE や Que2Code 等が GitHub に登録されています。

## 6. 現状の評価

ある時 Web 翻訳の精度が格段に向上したことに気づき驚いた記憶があります。その頃、海外向け製品のマニュアル校正で担当部分に間違いを見つけ、米国の編集者に「表 xx の yy 番目に列を追加し、その列見出しに…」と修正の指示をしたいのですが、そんな英語力はありません。社内に翻訳サイトがあったのですがインターネットに公開されている他の Web 翻訳と同様、英訳、和訳と往復させると単語の数が減っていくような代物です。ところが半ば自棄で検索サイトの翻訳機能を試したところ、関係代名詞を使いこなした上に英訳和訳を往復させても 2 往復目以降は表現が変化しません。これがちょうど BERT が公開された頃の出来事で、BERT (技術要素でいうと他の言語モデルでも使われている Embedding、Attention、Transformer) により単語と熟語レベルの翻訳から句点で区切られた文という単位で翻訳できるように人知れず進化していたわけです。

事前学習済言語モデルは数値化したトークン<sup>10</sup>とトークンの関連 (重み) 付けで保存された情報で、コード生成を行う場合は既存のソースコードから事前学習した言語モデルをトークン化した自然言語 (ソースコメント) で一致度を基準に検索します。これにより自然言語の文からソースコードの文や式に変換できるのですが、仕組みの観点から以下の点を考慮する必要があります。

- ① ソースコード 1 文を作るのにコメント 1 文を必要とするのであれば、コードを書いた方が早い
- ② ローカルで作成したクラスやメソッドは事前学習に含まれないので、初期は候補に含まれない

---

<sup>7</sup> <https://huggingface.co/models>

<sup>8</sup> HF Code Autocomplete

<https://marketplace.visualstudio.com/items?itemName=HuggingFace.huggingface-vscode>

<sup>9</sup> Stack Overflow <https://stackoverflow.com/search?q=stack+overflow+search+engine>

<sup>10</sup> トークンは英語系では単語 (空白等の区切り文字)、日本語の場合はもう少し細くなるようです…詳細は「トークン 日本語 Embedding」で検索

## コード自動生成 (CodeGeeX 他)

①のコメントと生成されるコードの関係は、学習が進むと 4.3.1 ②のようにコメントからメソッドや for each ブロック程度は生成できるようです。課題は日本語の場合はどう書くか、for each 以外の処理ブロックは生成させられないのか等があります。

②については、ローカルのクラスやメソッドを使っているコードを学習 (言語モデルへの登録) すれば提案候補に上がるようになると考えられます。また、設計工程を VSCode 等の開発環境を使って行えるようにすれば、設計工程の用語が言語モデルに取り込まれて設計/仕様から直接コードを生成することが可能になるかもしれません。

### 7. 今後

#### <背景>

現状 Ruby 以外の言語では日本人は言語仕様に関わることができておらず、スタックオーバーフロー等の国際的なコミュニティへの書き込みも見かけません。言語仕様やフレームワークの策定に関われないので頻繁に行われるバージョンアップへの追従が困難になります。新機能を使う予定がなくとも OS/ミドルウェア、ライブラリのセキュリティ対応や互換維持のためのバージョンアップが必要になります。これに伴って API が変わったり、廃止になったり、記法の変更になったりで古いアプリの改修が必要になり、現状維持のためだけの後ろ向きな仕事が発生します。

#### <その他の考慮点>

日本語を使ってコード再生成ができればバージョンの問題を軽減できる可能性があります。但し、古いアプリをバージョンアップさせるケースではコメントと生成したコードが混在する現状のツールでは困難で、仕様からコードを生成させるような今よりも高度な機能が必要になります。

#### <おすすめ>

事前学習済言語モデルを使ったコード生成ツールは定番 (慣用) コードの入力を省いてくれます。これだけでも有用性が高いですが、オープンソースの事前学習済言語モデルを取得して自前の開発物で強化学習を行えば、自組織が開発したクラスやメソッド、コードパターンを熟知した言語モデルが作れ、セキュリティも維持できます。

現状のツールが提示する内容は完全に満足できるものではなくスピードも遅いですが、ツールの開発は活発に進められており、いつか BERT が登場したときのように、気が付いたら景色が変わっていたということになるかもしれません。組織としても個人としても試してみる価値はあります。

※ツール及び環境の進展が早く、資料の内容を最新の情報に更新・維持し続けるのは困難なので、新しい追記が必要な情報は別の資料として纏めます。

この資料は 2023 年 4 月末時点の情報で作っています。

以上