



# TypeBERT (ESEC/FSE 2021) と ManyTypes4TypeScript (MSR 2022) の背景レビュー

オプショナル型言語TypeScriptの型推論研究として、JesseらはBERTスタイルの事前学習モデル「TypeBERT」を提案し<sup>1</sup>、Mirらは大規模データセット「ManyTypes4TypeScript」を構築・公開している<sup>2</sup>。以下、それぞれの研究におけるモデル規模、学習データ規模、型語彙制限、評価手法などを整理する。

## 大規模データ vs モデルサイズ

- **TypeBERT:** BERTに類する大規模モデル（エンコーダ層24層・隠れ層次元1024、総パラメータ約3.4億）を用い、JavaScript約2.5万プロジェクトで事前学習、TypeScript約2万プロジェクトで微調整している<sup>3</sup><sup>4</sup>。これは従来の数百プロジェクト規模のデータセットに比べ格段に大きく、Jesseら自身も「十分なデータがあれば単純なモデルでも最先端を超える性能を示せる」と述べている<sup>1</sup><sup>5</sup>。実際、TypeBERTは上位100型においてTop1精度約89.5%とLambdaNet（66.9%）を大幅に上回り、全体でも71.12%対64.2%と優位に立った<sup>6</sup><sup>7</sup>。一方で、この規模のモデルは学習コストも大きく、事前学習には6 GPU×160時間を要したと報告されている<sup>8</sup>。推論時は高速で、バッチ(64×256長)当たり1.28秒（1シーケンス約0.02秒）で処理できるとされる<sup>9</sup>。
- **ManyTypes4TypeScript:** TypeScript向けに新規構築された大規模データセット。13,953プロジェクト・539,571ファイルから約900万件の型注釈を集めており、Python向けデータセットの10倍に相当する規模である<sup>2</sup>。このデータセットは主に学習用に用いられ、HuggingFace上で提供されている。TypeBERTのような巨大モデルだけでなく、CodeBERTやGraphCodeBERT（110–125M程度）など比較的軽量なモデルでも学習・評価できるよう設計されている。
- **軽量モデルの優位／LLMの制約:** TypeBERTの例では高精度が得られたが、モデルサイズ3.4億では依然LLM（例えばGPT-3の1750億）には遠く及ばない。むしろBERT系モデルでは、モデル容量を抑えつつ大規模データで学習するアプローチが注目される。実際、TypeBERTは膨大なデータを消費する一方で、依存関係グラフ不要で高速推論が可能な点も強みである<sup>9</sup>。例えばTypeBERTはGPU一枚で16,384個の型推論を0.02秒/件で行え、LambdaNetが依存グラフ構築に失敗するケースでも問題なく処理できると報告されている<sup>9</sup>。このように、大規模学習を活用しつつも軽量モデルで高精度・高速処理を実現する方向性が示唆されている。

## TypeBERTの型語彙制限と評価手法

- **型語彙の制限:** TypeBERTでは出力側に学習データで出現頻度上位40,000種の型のみを語彙として用い、それ以上の型はすべてUNK（未知）ラベルにまとめた<sup>10</sup>。これによりテストデータ中のUNK出現率は約8%に抑えられている（テスト注釈の92%以上をカバー）<sup>10</sup><sup>11</sup>。ただし語彙を閉じるため、学習データに現れない型（珍しいユーザ定義型など）はUNKとなり、当該位置は誤答とカウントされる。
- **評価指標:** 予測精度の評価にはTop-1/Top-5分類精度を用い、「上位100型」（頻出型）「ユーザ定義型」「その他型」といったカテゴリに分けて測定している<sup>12</sup><sup>13</sup>。例えばTypeBERTは上位100型でTop1≈89.5%、Top5≈98.5%と非常に高い精度を示し<sup>14</sup>、開発者への推薦候補として実用的な性能

を備えた。一方で、モデルの予測が実際にプログラムで型検査に通るかどうか（型整合性）は評価しておらず、あくまでデータセットの注釈との一致度で性能を判断している。

- **その他の特徴:** 学習・評価時には、情報量が低い「any」型は除外される<sup>10</sup> <sup>15</sup>。TypeBERTはユーザ定義型を含む広い型にも対応し、LambdaNetよりも高い全体Top1精度（71.12% vs 64.2%<sup>6</sup>）を達成した。これは、複雑な静的解析バイアスなしに文脈から型を推測するBERTモデルの表現力の高さを示唆する結果である<sup>16</sup> <sup>7</sup>。

## ManyTypes4TypeScriptのデータスキーム

- **データ収集・構造:** GitHub上のTypeScriptプロジェクトをGraphQLクエリで約29,500件収集し、依存パッケージをpnpmで解決してから各ファイルのASTを解析している<sup>17</sup>。AST走査により、人手で注釈された型とコンパイラが推論した型を両方抽出し、コードのトークン列と対応する型ラベルのペアを生成する<sup>17</sup>。データ重複検出ツールで18%の重複コードを排除した後、型注釈付きファイルは約53.96万件に整理された<sup>18</sup>。最終的に13,953プロジェクト、539,571ファイル、9百万件の型注釈からなるコーパスが得られている<sup>2</sup> <sup>19</sup>。

- **データ分割と形式:** コーパスはプロジェクト単位で学習/検証/テストに約80/10/10%の割合で分割される<sup>20</sup>。HuggingFace Datasets形式で提供され、各サンプルはコードの単語トークン列（JSONのtokensリスト）と対応する型ラベル列（labelsリスト）からなる<sup>21</sup>。この構造により、任意のサブワードトークナイザやTransformerモデルで容易に扱える。

- **型注釈の特徴:** 注釈位置を見ると、変数宣言が約380万件、関数のパラメタが約370万件で最も多い<sup>22</sup>（図2）。型分布では文字列型、any型、数値型が上位を占め、特にanyが非常に多い<sup>23</sup>（図3）。any型は有用性が低いため学習・評価から除外されている<sup>24</sup>。推論可能な型割合（型カバレッジ）について、ManyTypes4TSでは「全ての箇所に注釈または推論型があれば静的型付き言語相当になる」と説明しており、現状未注釈部分を確率的手法で埋めることで型カバレッジを飛躍的に高められる可能性がある<sup>25</sup>。

- **型語彙と推論型割合:** 出力語彙は頻度上位50,000種に固定し、それ以上の型はUNKに置換する<sup>24</sup> <sup>26</sup>。この設定で型出現頻度の約94.1%をカバーする（残り約5.9%はUNK扱い）<sup>26</sup>。また、コーパス全体の型のうち約57%がコンパイラ推論型、43%が人手注釈型であり、特にany型では人手注釈が20%程度に留まる<sup>27</sup>。これらのデータ構造を用いることで、様々なコード特化モデルが統一的に訓練・評価可能となっている。

以上のように、TypeBERTは大規模データと高キャパシティモデルで高い分類精度を達成した例であり、ManyTypes4TypeScriptは大規模学習用コーパスの標準化を狙った取り組みである。両者を通じて、先行研究では数万～数十万規模のコードデータと数億規模のパラメータを組み合わせ、型語彙を上位数万種に限定して評価する手法がとられてきたことが明らかである<sup>3</sup> <sup>26</sup>。これら知見は、BERT系のより軽量なモデル設計においても、**モデルサイズ・データ量・語彙範囲**のトレードオフを考慮する上で重要な背景情報となる。

**参考文献:** TypeBERT <sup>1</sup> <sup>3</sup> <sup>12</sup> <sup>16</sup>、ManyTypes4TypeScript <sup>2</sup> <sup>17</sup> <sup>26</sup> <sup>28</sup> 他。

---

<sup>1</sup> <sup>3</sup> <sup>4</sup> <sup>5</sup> <sup>6</sup> <sup>7</sup> <sup>8</sup> <sup>9</sup> <sup>10</sup> <sup>11</sup> <sup>12</sup> <sup>13</sup> <sup>14</sup> <sup>15</sup> <sup>16</sup> Learning Type Annotation: Is Big Data Enough?  
[https://www.cs.ucdavis.edu/~devanbu/typebert\\_esec\\_fse\\_.pdf](https://www.cs.ucdavis.edu/~devanbu/typebert_esec_fse_.pdf)

<sup>2</sup> <sup>17</sup> <sup>18</sup> <sup>19</sup> <sup>20</sup> <sup>21</sup> <sup>22</sup> <sup>23</sup> <sup>24</sup> <sup>25</sup> <sup>26</sup> <sup>27</sup> <sup>28</sup> ManyTypes4TypeScript: A Comprehensive TypeScript Dataset for Sequence-Based Type Inference  
<https://www.kevinrjesse.com/pdfs/ManyTypes4TypeScript.pdf>