

# Phương Pháp Tạo Dữ Liệu cho Collaborative Fast & Slow Thinking Systems

## 1. Tổng quan về phương pháp tạo dữ liệu

### 1.1 Mục tiêu

Phương pháp tạo dữ liệu này nhằm xây dựng một bộ dữ liệu toàn diện để huấn luyện mô hình có khả năng tự động nhận biết và chuyển đổi giữa fast thinking và slow thinking. Bộ dữ liệu cần đảm bảo: - Đa dạng về loại nhiệm vụ và độ phức tạp - Phân biệt rõ ràng giữa fast thinking và slow thinking - Bao gồm cả quá trình tư duy và kết quả cuối cùng - Có thể mở rộng và cập nhật dễ dàng

### 1.2 Cấu trúc dữ liệu cải tiến

Mỗi mẫu trong bộ dữ liệu sẽ có cấu trúc sau:

```
{
  "id": "unique_id",
  "task": {
    "instruction": "Nhiệm vụ cần giải quyết",
    "input": "Dữ liệu đầu vào (nếu có)",
    "context": "Ngữ cảnh bổ sung (nếu có)"
  },
  "analysis": {
    "complexity_score": 0.75,
    "task_type": "reasoning | creative | qa | decision | ...",
    "constraints": ["Ràng buộc 1", "Ràng buộc 2", ...],
    "recommended_strategy": "fast_only | slow_only | fast_then_slow | parallel | iterative"
  },
  "thinking_processes": {
    "fast_thinking": {
      "process": "Quá trình tư duy nhanh, trực giác",
      "simplified_task": "Nhiệm vụ đã được đơn giản hóa",
      "intermediate_result": "Kết quả trung gian từ fast thinking"
    },
    "slow_thinking": {
      "decomposition": ["Bước 1", "Bước 2", ...],
      "reasoning": ["Suy luận cho bước 1", "Suy luận cho bước 2", ...],
```

```

    "verification": ["Xác thực kết quả bước 1", "Xác thực kết quả bước 2", ...],
    "intermediate_result": "Kết quả trung gian từ slow thinking"
  },
},
"integration": {
  "process": "Quá trình tích hợp kết quả từ fast thinking và slow thinking",
  "inspection": "Quá trình kiểm tra tính chính xác của kết quả",
  "corrections": ["Sửa lỗi 1", "Sửa lỗi 2", ...]
},
"output": {
  "final_answer": "Câu trả lời cuối cùng",
  "explanation": "Giải thích cho câu trả lời"
},
"metadata": {
  "domain": "math|logic|language|science|...",
  "difficulty": "easy|medium|hard|expert",
  "source": "synthetic|human|augmented|benchmark",
  "quality_score": 0.95,
  "tags": ["Tag 1", "Tag 2", ...]
}
}

```

## 2. Nguồn dữ liệu

### 2.1 Dữ liệu từ các benchmark hiện có

Kết hợp dữ liệu từ các benchmark hiện có và chuyển đổi sang định dạng mới:

#### 2.1.1 Dữ liệu cho nhiệm vụ suy luận

- **GSM8K, MATH:** Bộ dữ liệu toán học với các bài toán đòi hỏi suy luận từng bước
- **LogiQA, CLUTRR:** Bộ dữ liệu suy luận logic
- **BIG-Bench:** Các nhiệm vụ đa dạng với độ phức tạp khác nhau
- **MMLU:** Bộ dữ liệu kiến thức đa lĩnh vực
- **GPQA:** Bộ dữ liệu câu hỏi khoa học phức tạp

#### 2.1.2 Dữ liệu cho nhiệm vụ sáng tạo

- **CommonGen, CommonGen-Hard:** Bộ dữ liệu tạo câu chuyện có ràng buộc
- **WritingPrompts:** Bộ dữ liệu tạo văn bản sáng tạo
- **StoryGen:** Bộ dữ liệu tạo câu chuyện dài

#### 2.1.3 Dữ liệu cho nhiệm vụ lập trình

- **HumanEval, MBPP:** Bộ dữ liệu lập trình với các bài toán đòi hỏi suy luận
- **RustBrain Dataset:** Bộ dữ liệu sửa lỗi trong mã nguồn Rust

#### 2.1.4 Dữ liệu cho nhiệm vụ trả lời câu hỏi

- **LongBench**: Bộ dữ liệu trả lời câu hỏi dựa trên nội dung dài
- **ASAP**: Bộ dữ liệu đánh giá câu trả lời

### 2.2 Dữ liệu tổng hợp mới

Tạo dữ liệu mới để bổ sung cho các benchmark hiện có:

#### 2.2.1 Dữ liệu tạo ra bởi LLMs lớn hơn

- Sử dụng các mô hình như GPT-4, Claude 3, Gemini, v.v. để tạo dữ liệu
- Áp dụng kỹ thuật "imitate, explore, self-improve" để tạo dữ liệu chất lượng cao
- Tạo dữ liệu cho cả fast thinking và slow thinking

#### 2.2.2 Dữ liệu tạo ra bởi con người

- Thu thập dữ liệu từ các chuyên gia trong các lĩnh vực khác nhau
- Ghi lại quá trình tư duy của con người khi giải quyết các vấn đề phức tạp
- Phân loại và gán nhãn dữ liệu theo cấu trúc mới

#### 2.2.3 Dữ liệu tăng cường từ các nguồn có sẵn

- Mở rộng dữ liệu từ các benchmark hiện có
- Thêm các bước tư duy chi tiết và quá trình tích hợp
- Gán nhãn độ phức tạp và chiến lược tư duy phù hợp

## 3. Quy trình tạo dữ liệu

### 3.1 Quy trình tổng thể

Quy trình tạo dữ liệu bao gồm bốn giai đoạn chính:

1. **Thu thập và phân loại dữ liệu ban đầu**
2. **Tạo dữ liệu cho các chiến lược tư duy khác nhau**
3. **Tăng cường và đa dạng hóa dữ liệu**
4. **Kiểm tra chất lượng và lọc dữ liệu**

### 3.2 Thu thập và phân loại dữ liệu ban đầu

#### 3.2.1 Thu thập dữ liệu

- Thu thập dữ liệu từ các benchmark hiện có
- Thu thập dữ liệu từ các nguồn mở khác

- Tạo dữ liệu mới bằng cách sử dụng LLMs lớn hơn

### **3.2.2 Phân loại dữ liệu**

- Phân loại dữ liệu theo loại nhiệm vụ
- Phân loại dữ liệu theo độ phức tạp
- Phân loại dữ liệu theo lĩnh vực

### **3.2.3 Gán nhãn dữ liệu**

- Gán điểm số độ phức tạp (0-1)
- Gán nhãn loại nhiệm vụ
- Gán nhãn chiến lược tư duy phù hợp

## **3.3 Tạo dữ liệu cho các chiến lược tư duy khác nhau**

### **3.3.1 Tạo dữ liệu cho chiến lược Fast-Only**

- Xác định các nhiệm vụ phù hợp với fast thinking
- Tạo quá trình tư duy nhanh và câu trả lời cuối cùng
- Đảm bảo độ đa dạng của dữ liệu

### **3.3.2 Tạo dữ liệu cho chiến lược Slow-Only**

- Xác định các nhiệm vụ phù hợp với slow thinking
- Tạo quá trình tư duy chậm với các bước phân rã, suy luận và xác thực
- Đảm bảo độ đa dạng của dữ liệu

### **3.3.3 Tạo dữ liệu cho chiến lược Fast-then-Slow**

- Xác định các nhiệm vụ phù hợp với cả hai loại tư duy
- Tạo quá trình tư duy nhanh, sau đó là tư duy chậm
- Tạo quá trình tích hợp kết quả từ cả hai loại tư duy

### **3.3.4 Tạo dữ liệu cho chiến lược Parallel**

- Xác định các nhiệm vụ có thể giải quyết bằng nhiều cách tiếp cận
- Tạo quá trình tư duy nhanh và tư duy chậm song song
- Tạo quá trình tích hợp kết quả từ cả hai quá trình

### **3.3.5 Tạo dữ liệu cho chiến lược Iterative**

- Xác định các nhiệm vụ phức tạp, đòi hỏi nhiều vòng lặp
- Tạo quá trình tư duy lặp đi lặp lại giữa fast thinking và slow thinking
- Tạo quá trình tích hợp và cải thiện kết quả qua các vòng lặp

## 3.4 Tăng cường và đa dạng hóa dữ liệu

### 3.4.1 Biến đổi dữ liệu

- Thay đổi ngữ cảnh, độ phức tạp, và định dạng
- Tạo các biến thể của cùng một nhiệm vụ với các mức độ phức tạp khác nhau
- Thay đổi cách diễn đạt của nhiệm vụ và câu trả lời

### 3.4.2 Tạo dữ liệu đối nghịch

- Tạo câu hỏi gây nhầm lẫn, câu hỏi bẫy
- Tạo câu hỏi phức tạp có thể giải quyết nhanh
- Tạo câu hỏi đơn giản nhưng đòi hỏi suy luận sâu

### 3.4.3 Tạo dữ liệu tự động

- Sử dụng LLMs lớn để tạo dữ liệu tự động
- Áp dụng kỹ thuật self-play và bootstrapping
- Sử dụng kỹ thuật data augmentation để tăng cường dữ liệu

## 3.5 Kiểm tra chất lượng và lọc dữ liệu

### 3.5.1 Kiểm tra tính chính xác

- Kiểm tra tính chính xác của câu trả lời cuối cùng
- Kiểm tra tính hợp lý của quá trình tư duy
- Kiểm tra tính nhất quán giữa quá trình tư duy và câu trả lời

### 3.5.2 Kiểm tra tính đa dạng

- Đảm bảo độ đa dạng về loại nhiệm vụ
- Đảm bảo độ đa dạng về độ phức tạp
- Đảm bảo độ đa dạng về lĩnh vực

### 3.5.3 Lọc dữ liệu

- Loại bỏ dữ liệu chất lượng thấp
- Loại bỏ dữ liệu trùng lặp
- Loại bỏ dữ liệu không phù hợp với cấu trúc mới

## 4. Công cụ và kỹ thuật tạo dữ liệu

### 4.1 Công cụ tạo dữ liệu tự động

#### 4.1.1 LLM-based Data Generator

- Sử dụng LLMs lớn để tạo dữ liệu tự động
- Thiết kế prompt template để tạo dữ liệu theo cấu trúc mới
- Tích hợp các kỹ thuật như few-shot learning và chain-of-thought prompting

#### 4.1.2 Data Augmentation Pipeline

- Xây dựng pipeline tăng cường dữ liệu
- Tích hợp các kỹ thuật biến đổi dữ liệu
- Tích hợp các kỹ thuật tạo dữ liệu đối nghịch

#### 4.1.3 Quality Control System

- Xây dựng hệ thống kiểm tra chất lượng dữ liệu
- Tích hợp các kỹ thuật lọc dữ liệu
- Tích hợp các kỹ thuật đánh giá dữ liệu

### 4.2 Kỹ thuật tạo dữ liệu nâng cao

#### 4.2.1 Imitation Learning

- Bắt chước quá trình tư duy của con người
- Bắt chước quá trình tư duy của LLMs lớn hơn
- Tạo dữ liệu theo định dạng chuẩn

#### 4.2.2 Exploration

- Khám phá không gian nhiệm vụ để tìm các nhiệm vụ thách thức
- Tạo nhiều giải pháp cho cùng một nhiệm vụ
- Cải thiện giải pháp dựa trên phản hồi

#### 4.2.3 Self-Improvement

- Sử dụng dữ liệu chất lượng cao để cải thiện quá trình tạo dữ liệu
- Áp dụng các kỹ thuật như supervised fine-tuning (SFT) và direct preference optimization (DPO)
- Lọc dữ liệu chất lượng thấp dựa trên các chỉ số như độ dài và perplexity

## 5. Phân phối và cân bằng dữ liệu

### 5.1 Phân phối dữ liệu theo loại nhiệm vụ

- **Suy luận:** 30% tổng số dữ liệu
- **Sáng tạo:** 20% tổng số dữ liệu
- **Trả lời câu hỏi:** 25% tổng số dữ liệu
- **Ra quyết định:** 15% tổng số dữ liệu
- **Khác:** 10% tổng số dữ liệu

### 5.2 Phân phối dữ liệu theo độ phức tạp

- **Đơn giản (0.0-0.3):** 20% tổng số dữ liệu
- **Trung bình (0.3-0.6):** 40% tổng số dữ liệu
- **Phức tạp (0.6-0.8):** 30% tổng số dữ liệu
- **Rất phức tạp (0.8-1.0):** 10% tổng số dữ liệu

### 5.3 Phân phối dữ liệu theo chiến lược tư duy

- **Fast-Only:** 20% tổng số dữ liệu
- **Slow-Only:** 20% tổng số dữ liệu
- **Fast-then-Slow:** 30% tổng số dữ liệu
- **Parallel:** 15% tổng số dữ liệu
- **Iterative:** 15% tổng số dữ liệu

### 5.4 Cân bằng dữ liệu

- Đảm bảo sự cân bằng giữa các loại nhiệm vụ
- Đảm bảo sự cân bằng giữa các mức độ phức tạp
- Đảm bảo sự cân bằng giữa các chiến lược tư duy
- Đảm bảo sự cân bằng giữa các lĩnh vực

## 6. Kế hoạch triển khai

### 6.1 Giai đoạn 1: Chuẩn bị

- Thiết lập cơ sở hạ tầng cho việc tạo và lưu trữ dữ liệu
- Phát triển các công cụ tạo dữ liệu tự động
- Thiết lập quy trình kiểm tra chất lượng dữ liệu

## 6.2 Giai đoạn 2: Thu thập và chuyển đổi dữ liệu hiện có

- Thu thập dữ liệu từ các benchmark hiện có
- Chuyển đổi dữ liệu sang định dạng mới
- Gán nhãn dữ liệu theo cấu trúc mới

## 6.3 Giai đoạn 3: Tạo dữ liệu mới

- Tạo dữ liệu cho các chiến lược tư duy khác nhau
- Tăng cường và đa dạng hóa dữ liệu
- Kiểm tra chất lượng và lọc dữ liệu

## 6.4 Giai đoạn 4: Đánh giá và cải thiện

- Đánh giá chất lượng của bộ dữ liệu
- Xác định các khoảng trống trong dữ liệu
- Cải thiện quy trình tạo dữ liệu

## 6.5 Giai đoạn 5: Phát hành và duy trì

- Phát hành bộ dữ liệu cho cộng đồng
- Duy trì và cập nhật bộ dữ liệu
- Thu thập phản hồi và cải thiện bộ dữ liệu

# 7. Ví dụ minh họa

## 7.1 Ví dụ cho chiến lược Fast-Only

```
{
  "id": "fast_only_001",
  "task": {
    "instruction": "Tính tổng của 15 và 27.",
    "input": "",
    "context": ""
  },
  "analysis": {
    "complexity_score": 0.1,
    "task_type": "arithmetic",
    "constraints": [],
    "recommended_strategy": "fast_only"
  },
  "thinking_processes": {
    "fast_thinking": {
      "process": "Đây là một phép tính cộng đơn giản. Tôi sẽ cộng trực tiếp 15 và
```



```

27.",
  "simplified_task": "Tính 15 + 27",
  "intermediate_result": "42"
},
"slow_thinking": {
  "decomposition": [],
  "reasoning": [],
  "verification": [],
  "intermediate_result": ""
},
"integration": {
  "process": "Kết quả từ fast thinking là đủ chính xác.",
  "inspection": "Kiểm tra: 15 + 27 = 42. Đúng.",
  "corrections": []
},
"output": {
  "final_answer": "42",
  "explanation": "Tổng của 15 và 27 là 42."
},
"metadata": {
  "domain": "math",
  "difficulty": "easy",
  "source": "synthetic",
  "quality_score": 0.98,
  "tags": ["arithmetic", "addition", "fast_thinking"]
}
}

```

## 7.2 Ví dụ cho chiến lược Slow-Only

```

{
  "id": "slow_only_001",
  "task": {
    "instruction": "Giải phương trình bậc hai:  $2x^2 - 5x - 3 = 0$ ",
    "input": "",
    "context": ""
  },
  "analysis": {
    "complexity_score": 0.65,
    "task_type": "reasoning",
    "constraints": [],
    "recommended_strategy": "slow_only"
  },
  "thinking_processes": {
    "fast_thinking": {
      "process": "",
      "simplified_task": "",
      "intermediate_result": ""
    },
  },
}

```

```

"slow_thinking": {
  "decomposition": [
    "Bước 1: Xác định hệ số a, b, c trong phương trình  $ax^2 + bx + c = 0$ ",
    "Bước 2: Tính delta =  $b^2 - 4ac$ ",
    "Bước 3: Tính nghiệm  $x_1 = (-b + \sqrt{\text{delta}}) / (2a)$ ",
    "Bước 4: Tính nghiệm  $x_2 = (-b - \sqrt{\text{delta}}) / (2a)$ "
  ],
  "reasoning": [
    "Từ phương trình  $2x^2 - 5x - 3 = 0$ , ta có  $a = 2$ ,  $b = -5$ ,  $c = -3$ ",
    "Delta =  $b^2 - 4ac = (-5)^2 - 4 \times 2 \times (-3) = 25 + 24 = 49$ ",
    " $x_1 = (-b + \sqrt{\text{delta}}) / (2a) = (5 + 7) / 4 = 12 / 4 = 3$ ",
    " $x_2 = (-b - \sqrt{\text{delta}}) / (2a) = (5 - 7) / 4 = -2 / 4 = -0.5$ "
  ],
  "verification": [
    "Kiểm tra  $x_1 = 3$ :  $2(3)^2 - 5(3) - 3 = 2(9) - 15 - 3 = 18 - 15 - 3 = 0 \checkmark$ ",
    "Kiểm tra  $x_2 = -0.5$ :  $2(-0.5)^2 - 5(-0.5) - 3 = 2(0.25) + 2.5 - 3 = 0.5 + 2.5 - 3 = 0 \checkmark$ "
  ],
  "intermediate_result": " $x_1 = 3$ ,  $x_2 = -0.5$ "
}
},
"integration": {
  "process": "Kết quả từ slow thinking là đủ chính xác.",
  "inspection": "Đã kiểm tra cả hai nghiệm và cả hai đều thỏa mãn phương trình gốc.",
  "corrections": []
},
"output": {
  "final_answer": " $x_1 = 3$ ,  $x_2 = -0.5$ ",
  "explanation": "Phương trình  $2x^2 - 5x - 3 = 0$  có hai nghiệm là  $x_1 = 3$  và  $x_2 = -0.5$ ."
},
"metadata": {
  "domain": "math",
  "difficulty": "medium",
  "source": "synthetic",
  "quality_score": 0.95,
  "tags": ["algebra", "quadratic_equation", "slow_thinking"]
}
}

```

### 7.3 Ví dụ cho chiến lược Fast-then-Slow

```

{
  "id": "fast_then_slow_001",
  "task": {
    "instruction":
      "Viết một câu chuyện ngắn (khoảng 100 từ) về một chú mèo tên là Whiskers, bao gồm các từ: cửa sổ, mưa, và sách.",
    "input": "",
    "context": ""
  },
}

```

```
"analysis": {
  "complexity_score": 0.45,
  "task_type": "creative",
  "constraints": ["Bao gồm từ 'cửa sổ'", "Bao gồm từ 'mưa'", "Bao gồm từ 'sách'",
  "Khoảng 100 từ"],
  "recommended_strategy": "fast_then_slow"
},
"thinking_processes": {
  "fast_thinking": {
    "process": "Tôi cần viết một câu chuyện ngắn về một chú mèo tên là Whiskers. Tôi sẽ tạo một câu chuyện đơn giản trước, sau đó sẽ đảm bảo bao gồm các từ yêu cầu và độ dài phù hợp.",
    "simplified_task": "Viết một câu chuyện ngắn về một chú mèo tên là Whiskers",
    "intermediate_result": "Whiskers là một chú mèo thích ngồi bên cửa sổ. Một ngày mưa, Whiskers quan sát những giọt mưa rơi xuống. Chủ của Whiskers đang đọc sách gần đó. Whiskers nhảy vào lòng chủ và cùng tận hưởng buổi chiều yên bình."
  },
  "slow_thinking": {
    "decomposition": [
      "Bước 1: Kiểm tra xem câu chuyện đã bao gồm tất cả các từ yêu cầu chưa",
      "Bước 2: Kiểm tra độ dài của câu chuyện",
      "Bước 3: Cải thiện câu chuyện để làm cho nó hấp dẫn hơn",
      "Bước 4: Đảm bảo câu chuyện có cấu trúc tốt với đầu, giữa và cuối"
    ],
    "reasoning": [
      "Câu chuyện đã bao gồm các từ 'cửa sổ', 'mưa', và 'sách'",
      "Câu chuyện có khoảng 40 từ, cần thêm chi tiết để đạt khoảng 100 từ",
      "Câu chuyện khá đơn giản, cần thêm chi tiết về tính cách của Whiskers và môi trường xung quanh",
      "Câu chuyện cần có một kết thúc rõ ràng hơn"
    ],
    "verification": [
      "Đã bao gồm tất cả các từ yêu cầu",
      "Sau khi thêm chi tiết, câu chuyện có khoảng 100 từ",
      "Câu chuyện đã hấp dẫn hơn với chi tiết về tính cách của Whiskers",
      "Câu chuyện có cấu trúc tốt với đầu, giữa và cuối"
    ],
    "intermediate_result": "Whiskers, một chú mèo mướp tinh nghịch với đôi mắt xanh biếc, có thói quen ngồi hàng giờ bên cửa sổ để quan sát thế giới bên ngoài. Một buổi chiều mưa rào, những giọt nước nháy múa trên kính cửa sổ thu hút sự chú ý của Whiskers. Chú mèo đặt chân trước lên kính, cố gắng bắt những giọt mưa đang trượt xuống. Trong căn phòng ấm áp, chủ của Whiskers, một người phụ nữ trẻ, đang cuộn tròn trên ghế bành với cuốn sách yêu thích. Tiếng mưa rơi và tiếng lật trang sách tạo nên một bản nhạc nhẹ nhàng. Whiskers, mệt mỏi với trò chơi của mình, nhảy xuống khỏi bệ cửa sổ, đuổi người, và nhẹ nhàng nhảy vào lòng chủ. Cả hai cùng tận hưởng buổi chiều yên bình, một người với cuốn sách, một mèo với giấc ngủ ngon."
  }
},
"integration": {
  "process": "Kết hợp ý tưởng ban đầu từ fast thinking với phiên bản được mở
```

rộng và cải thiện từ slow thinking.",

**"inspection":** "Kiểm tra lại câu chuyện để đảm bảo nó đáp ứng tất cả các yêu cầu: bao gồm các từ yêu cầu, độ dài phù hợp, và có cấu trúc tốt.",

**"corrections":** ["Điều chỉnh một số từ để làm cho câu chuyện tự nhiên hơn"]

},

**"output":** {

**"final\_answer":** "Whiskers, một chú mèo mướp tinh nghịch với đôi mắt xanh biếc, có thói quen ngồi hàng giờ bên cửa sổ để quan sát thế giới bên ngoài. Một buổi chiều mưa rào, những giọt nước nhảy múa trên kính cửa sổ thu hút sự chú ý của Whiskers. Chú mèo đặt chân trước lên kính, cố gắng bắt những giọt mưa đang trượt xuống. Trong căn phòng ấm áp, chủ của Whiskers, một người phụ nữ trẻ, đang cuộn tròn trên ghế bành với cuốn sách yêu thích. Tiếng mưa rơi và tiếng lật trang sách tạo nên một bản nhạc nhẹ nhàng. Whiskers, mệt mỏi với trò chơi của mình, nhảy xuống khỏi bệ cửa sổ, đuổi người, và nhẹ nhàng nhảy vào lòng chủ. Cả hai cùng tận hưởng buổi chiều yên bình, một người với cuốn sách, một mèo với giấc ngủ ngon.",

**"explanation":**

"Câu chuyện ngắn về chú mèo Whiskers bao gồm các từ yêu cầu: 'cửa sổ', 'mưa', và 'sách'. Câu chuyện có khoảng 100 từ và có cấu trúc rõ ràng với đầu, giữa và cuối."

},

**"metadata":** {

**"domain":** "language",

**"difficulty":** "medium",

**"source":** "synthetic",

**"quality\_score":** 0.92,

**"tags":** ["creative\_writing", "story", "fast\_then\_slow"]

}

}

## 8. Kết luận

Phương pháp tạo dữ liệu này cung cấp một cách tiếp cận toàn diện để xây dựng bộ dữ liệu chất lượng cao cho việc huấn luyện mô hình có khả năng tự động nhận biết và chuyển đổi giữa fast thinking và slow thinking. Bằng cách kết hợp dữ liệu từ các benchmark hiện có với dữ liệu tổng hợp mới, phương pháp này đảm bảo độ đa dạng và chất lượng của bộ dữ liệu. Quy trình tạo dữ liệu chi tiết và các công cụ tạo dữ liệu tự động giúp triển khai phương pháp này một cách hiệu quả. Phân phối và cân bằng dữ liệu đảm bảo bộ dữ liệu phù hợp cho việc huấn luyện mô hình theo khung khái niệm cải tiến.