



浙江财经大学

本科生专业实践报告

题目：画像生成および処理システムの設計と実装

学生氏名： 吴国涛

学生番号： 220110900830

指導教員： 张睿峰

所属学院： 情報技術・人工知能学院

専攻名称： 人工智能

クラス： 二組

2025 年7 月

画像生成および処理システムの設計と実装

要旨: 視覚コンテンツの制作および処理の需要が増加し続ける中、画像の生成、編集、および理解技術は多くの応用シーンで重要な役割を果たしています。本文では、テキストからの画像生成（DashScope ImageSynthesisに基づく）、トライマップ不要の画像の切り抜き（MODNetに基づく）、テキスト認識・抽出、および画像のトリミングと解像度調整機能を一体化した画像処理システムを構築し、複雑なマルチモーダルタスクの一体化処理効率の向上を目指しました。まず、画像生成の段階では、DashScope SDKを採用して通義万相（Tongyi Wanxiang）シリーズの大規模モデルを呼び出し、高品質なテキスト生成画像の出力を実現しました。次に、画像の切り抜き機能は、軽量かつリアルタイムのトライマップ不要な人物切り抜きモデルであるMODNetに基づいており、前景と背景の高速な分離を実現すると同時に、詳細と意味的融合の高品質な効果を維持しています。そして、システムはOCR（光学文字認識）モジュールを統合し、画像内のテキスト情報を自動的に抽出し、情報処理能力を強化しました。最後に、画像のトリミングと解像度調整機能を組み込むことで、システムに柔軟な編集能力と多様な画像の出力能力を持たせました。全体として、本文のシステムはモジュール化された設計と機能統合を通じて、テキストプロンプトからの画像生成、画像分割、情報抽出から画像編集までのループを実現し、画像処理プロセスの自動化と対話効率を向上させました。

キーワード: DashScope ImageSynthesis, MODNet, OCR, 文生図（テキスト画像生成）

目次

要旨	I
1 はじめに	1
1.1 研究背景	1
1.2 国内外の研究現状	2
1.2.1 テキスト画像生成（文生図）の研究現状	2
1.2.2 テキスト画像生成（文生図）の研究現状	2
1.2.3 人物認識の研究現状	3
1.3 本文の内容	5
2 重要技術	6
2.1 DashScope	6
2.2 MODNet	6
2.3 EasyORC	7
3 要求分析とシステム概要	8
3.1 実現可能性分析	8
3.1.1 システムの実現可能性分析	8
3.1.2 時間的実現可能性分析	10
3.2 要求分析	10
4 システム実装	11
4.1 開発環境	11
4.2 主要モジュールの説明	11
4.2.1 登録およびログインページ	12
4.2.2 メインページ	13
4.2.3 文字抽出	14
4.2.4 画像切り抜き（Matting）	15
4.2.5 テキスト画像生成（文生図）	16
4.2.6 トリミングと解像度変更	17

5	テキスト認識、人物認識、テキスト画像生成	18
5.1	EasyORC	18
5.2	MODNet	18
5.3	DashScope	19
5.4	大規模モデルによるテキスト画像生成	20
6	評価	22
7	まとめと展望	23

1 はじめに

1.1 研究背景

視覚コンテンツの生成、編集、および理解は、現代のマルチメディアアプリケーションにおける核心的な需要となりつつあります。社会経済の持続的な発展とデジタル化の波に後押しされ、高品質な視覚体験に対する大衆の需要も日々高まっており、画像処理技術の絶え間ない進化と融合を促進しています。

テキスト画像生成（文生図）は、近年の大規模モデルの視覚能力における重要な表現形式です。アリババクラウド DashScope が提供する通義万相（Tongyi Wanxiang）シリーズモデル（wanx-v1、wanx2.1-t2i-turbo など）に代表されるように、このモデルは中英バイリンガル入力、多様な画風生成をサポートし、さらには参考画像に基づいたコンテンツやスタイルの転送も可能であり、生成品質が高く、呼び出しが便利であるという特徴を持っています。DashScope は Model-as-a-Service プラットフォームとして、開発者が複雑なインフラを構築することなく、SDK を通じてモデルを呼び出し画像を生成することを可能にします。

一方で、画像の切り抜き（Matting）、特に人物切り抜きにおける前景抽出と背景分離も、画像処理チェーンにおいて不可欠な要素です。MODNet は香港城市大学と商湯科技（SenseTime）によって提案された軽量かつリアルタイムの「トライマップ不要（Trimap-free）」な人物切り抜きモデルであり、RGB 入力のみで高精度な切り抜きを実現し、極めて高い実行効率を備えています。その構造設計には、意味予測、詳細予測、意味-詳細融合の 3 つの主要モジュールが含まれており、SOC 自己教師あり一貫性戦略を通じてモデルの堅牢性と汎化能力を向上させています。

画像・テキスト処理フローにおいて、テキスト認識（OCR）モジュールは画像中の文字情報を抽出するために使用され、シーン理解とコンテンツ管理をサポートします。また、画像のトリミングと解像度調整機能は、実際の応用におけるシステムの柔軟性とカスタマイズ能力をさらに強化し、様々な端末での表示や保存の需要に適応させます。

画像生成、切り抜き、文字認識、編集技術はそれぞれ実践において広く研究・応用されていますが、それらを一体化した多機能閉ループシステムに統合することは、依然として多くの課題に直面しています。これには、モジュール間のインターフェースとデータフローの統合の複雑さ、マルチタスク同時実行時に発生しうるパフォーマンスのボトルネック、各モジュール（特に OCR と生成モジュール）の生成結果における意味的一貫性のエラー、モジュールを跨ぐエラーの追跡とデバッグの難易度が高くシステムの安定性維持が困難であることなどが含まれますが、これらに限定されません。

したがって、テキスト画像生成、トライマップ不要の切り抜き、画像内コンテンツ認識、編集機能等を一体化したモジュール式システムを構築することは、視覚処理

タスクの自動化と効率を向上させるだけでなく、強力な対話性を備えており、著しい応用価値があります。

1.2 国内外の研究現状

1.2.1 テキスト画像生成の研究現状

テキスト生成画像 (Text-to-Image, T2I) は、自然言語の記述を自動的にリアルな画像に変換する重要な AI タスクであり、近年この分野は著しい進展を遂げています。現在の主な技術ルートには、条件付き敵対的生成ネットワーク (GAN)、自己回帰モデル、拡散モデル (Diffusion Models) などが含まれ、StackGAN、AttnGAN、DALL-E、Stable Diffusion、Imagen、RAPHAEL などの一連の象徴的なモデルと技術が生まれています。その中で、GAN モデルはエンドツーエンドの学習を通じてテキストと画像の直接的なマッピングを実現します。例えば、StackGAN は 2 段階の生成構造を利用してテキストのスケッチを段階的に高品質な画像へと精細化します。一方、Imagen や RAPHAEL などの拡散モデルは、ノイズの追加と除去を段階的に行うことで高忠実度な画像を生成し、Imagen は 2025 年に発表されたバージョン Imagen 4 まで発展しており、その画像の細部とリアリズムは大幅に向上しています。総説研究も示しているように、拡散モデルは現在最も将来性のある技術ルートの一つであり、多くの研究がその生成性能の向上、マルチモーダルな一貫性、画像制御能力、安全性などに焦点を当てています。さらに、データ不足の問題を解決するために、自己教師あり GAN、条件付き拡張、テキスト意味一貫性の最適化などの手法も広く研究されており、SS-TiGAN のように自己教師ありメカニズムと特徴マッチング手法を提案してテキストと画像の一貫性および視覚品質を改善する例もあります。

1.2.2 文字認識の研究現状

文字認識 (OCR) の研究は全体として、初期のテンプレートや特定のフォントに依存するものから、現代のディープラーニングを採用する方法へと進化しており、現在の主な傾向は以下の通りです。従来の OCR (Tesseract など) はテンプレート

認識や印刷フォント基準に基づいており、「きれいな」ドキュメントに対しては安定して信頼性が高く、レイアウト分析やhOCR出力をサポートするため、オフライン処理や公益的なシーンに適しています。近年ではディープラーニングOCRモデルが広く応用されています。例えば、EasyOCR（CRAFTテキスト検出とCRNN認識の結合）は多言語サポートと使いやすいインターフェースを提供し、シーン内の文字に対して良好なパフォーマンスを示します。PaddleOCRは中国語市場で特に際立っており、高効率な検出、認識、構造化出力などの機能をサポートし、正確さと速度のバランスをとっています。MMOCRやdocTRなどの研究ツールボックスは、完全なOCRパイプライン、柔軟なモデル選択、現代的なアーキテクチャ（Transformerサポートなど）を提供し、開発者が複雑なドキュメント認識タスクに適応するのに便利です。歴史的文書や多言語入力の処理に関しては、Kraken（OCRopusの後継）やOCRopus自体が古籍や非ラテン文字の認識によく使用されます。CalamariはFraktur（亀甲文字）などのフォントに対しても顕著な認識性能を備えています。要約すると、Tesseractは広範な互換性と成熟度を持ち、EasyOCRとPaddleOCRは実用的なシーンでの言語サポートと統合の利便性を重視しており、MMOCR/docTRなどのツールは実験や複雑な構造化アプリケーションに適しています。

1.2.3 人像识别研究现状

近年、トライマップ不要（Trimap-free：追加の注釈不要）の人物切り抜き技術が著しく進展しており、ビデオ通話、ライブ配信、AR/VRなどのシーンにおけるリアルタイム性、軽量化、細部再現への高い要求に応えています。MODNet（Matting Objective Decomposition Network）はこの分野における画期的な成果となりました。単一のRGB画像入力を通じて、目的分解方式により意味推定、境界詳細、融合を同時に最適化し、e-ASPPモジュールを利用してマルチスケールの情報を融合し、自己教師ありSOC（サブ目標一貫性）戦略を用いて実シーンへの汎化能力を向上させています。Adobe Matting DatasetおよびPPM-100データセットにおいて、それまでのトライマップ不要の手法よりも大幅に優れていました。その軽量の構造は7MBまで圧縮可能で、通常のPCや携帯電話での2K解像度画像の高速処理をサポートします。

その後の研究は、主に3つの方向へ展開しています。第一に精細レベルの向上です。例えばPP-Mattingは高解像度詳細ブランチ（HRDB）と意味コンテキストブランチ（SCB）を採用し、Trimapに依存せずにalpha matteの精度と詳細表現を向上させました。第二に高解像度性能の最適化です。Yatao Zhongらが提案した2段階アーキテクチャは、ViTを導入して粗い推定を行い、CRA（Cross Region Attention）モジュールを通じてエッジを精修し、HD/4K近いビデオでほぼリアルタイムの推論を実現しつつ、FLOPSを大幅に（元のモデルの約1/20に）削減しました。第三に実シーンでの汎化強化です。MFC-Netは多特徴融合メカニズムを利用し、史上最大規模のリアル・トライマップ不要切り抜きデータセットReal-19kを構築し、実画像でのモデルのパフォーマンスを著しく向上させました。

同時に、研究者はトライマップ不要の切り抜きをより複雑または特殊なシーンに応用することも進めています。例えば、2025年に提案されたFlash-Priorsモデルは、フラッシュあり/なしの画像ペアを入力とし、transformerモジュールと境界精修ネットワークを組み合わせ、動的背景下での切り抜きの堅牢性と精度を向上させました。また、SFMattingNetはこのパラダイムを煙や炎などの非剛体ターゲットに拡張し、SFMatting-800データセットを構築し、このタスクにおいてMODNetと比較して平均誤差を約12.65%低減させました。

全体として、トライマップ不要のリアルタイム軽量人物切り抜きの発展は以下の傾向を示しています：1. エンドツーエンドの軽量アーキテクチャによるリアルタイム推論の加速、2. 粗・精の結合とマルチタスク協調による精度の向上、3.

Transformerなどの現代的構造の導入による高解像度処理能力の最適化、4. 実シーンデータセットの構築による汎化能力の強化、5. より多くの複雑なシーンやクラス横断的なターゲットへの応用拡大。

1.3 本文内容

本文の主な貢献は以下のように要約されます。

- 1) 軽量かつ実用的な画像処理システムを設計・実装しました。ログイン/登録、文生図（テキスト生成画像）、人物切り抜き、文字認識（OCR）、解像度変更とトリミングなどの機能を同一プラットフォームに統合し、ユーザー認証から画像生成、編集、エクスポートまでの完全なワークフローを確立しました。このシステムは統一されたインターフェースとフレンドリーなフロントエンドのインタラクションを通じて、ユーザーの「テキスト入力/画像アップロード」から「使用可能な画像の産出」までの操作手順を大幅に簡素化し、ユーザー体験と使用効率を著しく向上させました。
- 2) 重要な画像処理モジュールにおいてエンジニアリング最適化を行いました。人物切り抜きと文字認識のために安定した推論パイプラインと後処理フロー（透明背景出力、多形式エクスポート、自動トリミングに対応）を構築し、モデルの軽量化と推論最適化、入力の前処理とキャッシュ戦略などの措置を通じて、精度を保証しつつ応答遅延を低減しました。同時に、基本的な認証と権限管理を実現し、ユーザーがアップロードしたデータと結果のプライバシー保護を保障しました。
- 3) 文生図と編集機能のために高効率な指令/テンプレートメカニズムとモジュール化アーキテクチャを設計しました。中国語の意味指令の解析とテンプレート化されたプロンプトをサポートし、ユーザーの自然言語記述をより安定した生成入力に変換できるようにしました。また、OCRと切り抜き結果を下流情報として自動補正やインテリジェントな後処理（例えば、認識した文字に基づき置換レイヤーを自動生成する、切り抜き結果に基づき背景を自動合成するなど）に利用し、生成コンテンツの関連性と編集可能性を向上させました。システムはモジュール化され拡張可能な設計を採用しており、将来的にさらなるモデルや能力を接続することを容易にし、若

干の定量的指標と小規模なユーザーテストを通じて、システムの正確性、応答速度、使いやすさを検証しました。

2 重要技術

機能が充実したシステムを開発するには、今日の多数のオープンソースのコンピュータビジョンおよびディープラーニング技術の支援が欠かせません。本文のシステムは、Dashscope ImageSynthesisを用いて画像の合成と強化を実現し、MODNetを利用して正確な前景切り抜きと透明度推定を実行し、EasyOCRの助けを借りて画像内のテキスト検出と認識を完了します。システムの可用性とデプロイ可能性を保証するため、プロジェクトの開発過程では、一般的なフロントエンド・バックエンドフレームワークとコンテナ化手段を組み合わせ、開発環境の統合を実現しました。本章では、本システムの開発に関わる各重要技術について詳細に紹介します。

2.1 DashScope ImageSynthesis

DashScopeはアリババクラウドが提供する「Model-as-a-Service」(MaaS)フレームワークであり、AIモデルのホスティングと呼び出しに特化しています。DashScope ImageSynthesisを利用することで、開発者はクラウド上の事前学習済みテキスト画像生成モデルを簡単に呼び出すことができ、GPUや複雑なインフラを自前で構築する必要がありません。その核心的な利点は以下の通りです。

インフラの分離: DashScopeはAPI接続方式を提供し、底層のモデルデプロイやハードウェア要件を隠蔽するため、開発者はビジネスロジックに集中できます。

多言語SDKサポート: 現在、PythonとJavaのSDKを提供しており、一般的なバックエンドサービスへの統合が容易です。

テキストプロンプト駆動の画像生成: `ImageSynthesis.call(model='qwen-v1', prompt=...)` のような方式で画像を構築し、URLまたは画像データを出力します。

クラウドプラットフォームに基づく課金： 回数課金方式（例：Qwen-Imageモデルは1枚あたり約\$0.035）を採用しており、無料枠や同時実行タスク制御もサポートしています。

2.2 MODNet

MODNet (Matting Objective Decomposition Network) は、軽量なリアルタイム人物切り抜きモデルであり、以下の重要な特徴を持っています。

- ・ **Trimap不要、高いリアルタイム性：** 単一のRGB画像のみで分割切り抜きが可能で、前景・背景領域に対する手動注釈への依存を回避します。
- ・ **高性能：** リアルタイムビデオや高解像度画像処理に適しています。
- ・ **軽量な3ブランチ構造：** 低解像度ブランチ（高速な意味推定用）、高解像度ブランチ（正確な境界検出）、融合ブランチ（情報を統合し最終的なalpha mattesを生成）。
- ・ **重要技術：** e-ASPP (Efficient Atrous Spatial Pyramid Pooling) モジュールを導入して特徴抽出効率を向上させ、SOC (Self-supervised Objectives Consistency) 戦略を通じて実シーンへのモデルの適応性を強化しています。

2.3 EasyOCR

EasyOCRはオープンソースのOCR（光学文字認識）ツールで、使いやすく、多言語をサポートし、迅速な統合に適しているという特徴があります。そのハイライトは以下の通りです。

- ・ **サポート言語が広範：** 中国語、英語、アラビア語、インド系文字（デーヴァナーガリーなど）を含む80以上の言語と多数の筆記スクリプトをサポートしています。

- ・ **即座に使用可能 (Plug-and-Play) :** `pip install easyocr` でインストール後、モジュールをインポートし、`reader = easyocr.Reader(['ch_sim'], 'en'])` でモデルをロードし、`reader.readtext(...)` を呼び出すだけで認識可能です。
- ・ **認識が正確で堅牢:** 内部でResNet、LSTM、CTCデコードメカニズムを組み合わせており、自然なシーンやドキュメント画像内の複雑な文字レイアウトに対応できます。
- ・ **設定が柔軟:** `readtext` メソッドはdecoder、batch_size、テキスト検出閾値、許容範囲 (allowlist)、回転検出などのパラメータ設定をサポートし、異なる応用シーンの需要を満たします。

3 要求分析とシステム概要

3.1 実現可能性分析

3.1.1 システムの実現可能性分析

本プロジェクトはWebベースの軽量画像処理サイトであり、Dashscope ImageSynthesis（画像合成/強化）、MODNet（高品質前景切り抜き/透明マスク生成）、EasyOCR（文字検出と認識）の3つの技術を核心機能として展開します。技術的実現の観点から見て、本システムは以下の理由により良好な実現可能性を持っています。

技術ソリューションが成熟しており相互補完的である：

- Dashscope ImageSynthesisは、画像コンテンツの生成/補完、スタイル転送やデータ拡張に使用でき、ユーザーに自動合成、背景置換、画像修復の能力を提供します。
- MODNetは前景切り抜き（alpha matte生成）に特化しており、高品質な透明前景レイヤーを出力でき、精細な背景除去や置換シーンに適しています。
- EasyOCRは、すぐに使える多言語文字検出・認識能力を提供し、画像からのテキスト情報抽出、編集可能な文字オーバーレイ、または文字領域ごとの局所処理に使用できます。この3者は機能的に補完し合っています。まずMODNetで前景/マスクを抽出し、次にDashscopeで合成/修復または背景置換を行い、最後にEasyOCRで文字を抽出またはオーバーレイすることで、「インテリジェント切り抜き → スタイル/コンテンツ合成 → テキスト認識編集」という完全なパイプラインを実現します。

実装とデプロイの技術スタック：

バックエンドにはPython（ディープラーニングモデルの統合、PyTorch/ONNXランタイムの呼び出し、またはFlask/FastAPIを使用した推論APIの提供に便利）を採用します。フロントエンドにはReact + SCSS（またはその他の現代的フレームワーク）を採用し、インタラクティブなUI（プレビュー、マスク微調整、テキストボックス編集、一括タスク管理など）を実現します。モデル推論はPyTorch、TorchScript、または

ONNXエクスポートを通じて行い、必要に応じてGPU（CUDA）サーバーやクラウド推論サービスに接続して性能要求を満たします。コンテナ化（Docker）と非同期キュー（Celery / Redisなど）は本番環境の拡張に使用可能です。

対話と可視化のサポート:

フロントエンドではリアルタイムプレビュー、マスク境界の微調整（ブラシ/消しゴム）、マルチビュー表示（元画像/前景/マスク/合成結果）、ズームや全画面表示などを実装でき、ユーザー体験を向上させるとともにモデル出力の人為的な補正を容易にします。

リソースとリスクが制御可能:

使用するモデルには成熟した事前学習済み重みと多数の参考実装が存在し、ゼロからの学習作業量を削減できます。主要なリスクは推論遅延とエッジケース（複雑な背景、極端な照明、低解像度テキスト）ですが、モデルの量子化、解像度適応、または「手動微調整」ツールの提供によって緩和可能です。

以上より、Dashscope ImageSynthesis、MODNet、EasyOCRを利用して構築する画像処理サイトは技術的に実現可能であり、3者の組み合わせによって切り抜きから合成、文字認識までの主要な機能シーンをカバーできます。

3.1.2 時間の実現可能性分析

本文システムの開発初期において、詳細なタスクスケジュール表を策定しました。タスクの重要なノード設定により、システム開発作業全体を効果的に細分化しました。問題定義の初期段階で、コミュニケーションを通じてシステムの期待目標を明確にし、それによってシステム機能を正確に定義し、方向性の逸脱を効果的に防ぎ、重複開発を回避しました。現在直面している課題を解決するために、広範な資料収集作業も実施し、大量の資料を収集することに成功しました。これらの前期準備は、後続の作業展開に強力なサポートを提供しました。システム開発過程では、反復開発（イテレーション）戦略を採用しました。システム開発プロセス全体を複数の段階的サイクルに分割することで、開発プロセスを持続的に監視し、潜在的な問題にタイムリーに対応することができ、システム開発が段階的な目標を予定通りに達成することを確実にしました。

3.2 要求分析

システム要求分析の面では、本システムはDashscope ImageSynthesisに基づく画像生成/編集、MODNetに基づく高品質前景切り抜きと透明背景出力、およびEasyOCRに基づく多言語テキスト検出と認識の3大機能モジュールをサポートする必要があります。エンドユーザー向けに直感的なWeb対話インターフェース（アップロード/レビュー/パラメータ調整/ダウンロード）と一括処理インターフェースを提供します。非機能要件としては、応答遅延が制御可能であること（中解像度画像1枚の処理目標 < 2 - 5秒、ハードウェアによる）、短時間のピークをサポートするための並行処理能力の拡張性、サーバー側でのGPUアクセラレーションのサポートおよびリソース不足時のCPUモードへのフォールバックが求められます。システムは入出力画像と認識テキストのプライバシーを保証し、転送層の暗号化、オプションのローカル一時保存および処理後の自動消去を採用すべきです。モデルのバージョン管理と設定可能なパラメータ（生成スタイル、切り抜き精度、OCR言語リストなど）、エラー処理と追跡可能なログを提供し、バックトラックと監視を容易にする必要があります。インターフェース層はREST/JSON仕様に準拠し、認証（API KeyまたはOAuth）

を提供し、フロントエンドは主流のブラウザと互換性があり、モバイル端末にも親和性がある必要があります。品質指標には、画像切り抜きのIoU/視覚的なシームレスさ、OCR認識の正確率、生成コンテンツの意味的一貫性が含まれ、将来的なモデルの置き換えや拡張（新規モデルや微調整）をサポートしつつ、コスト管理（GPU/帯域幅）と保守性（自動デプロイ、単体/結合テスト、ドキュメント）も考慮する必要があります。

4 システム実装

4.1 開発環境

1. システム環境: Windows 11
2. 開発ツール: Pycharm、Edge
3. データベース: SQLite
4. 言語: Python3.9
5. 前後端フレームワーク: Fastapi、Streamlit

4.2 主要モジュールの説明

本文のシステムインターフェースは白を基調色として採用しています。ユーザーが異なるシステム画面間を切り替えたり、必要なコンテンツを素早く見つけたりできるように、ページの左側にナビゲーションバーを設計しました。同時に、右側の領域はコンテンツ表示用に確保されており、操作画面に十分なスペースを確保し、ユーザーの多様な対話ニーズを満たしています。

4.2.1 登録およびログインページ

ユーザー登録ページでは、ユーザーは未登録のユーザー名を使用して登録を行うことができます。登録フローが完了すると、システムは自動的にユーザーをログインへ誘導します。登録に成功しなかった場合、システムはユーザーに登録失敗の具体的な原因を明確に提示します。

注册新账号

用户名（注册）

密码

重复密码

注册并登录

返回登录

ユーザーログインページでは、ユーザーは自分のアカウント情報を入力してログインできます。ユーザーが提出した情報がシステムによって検証され、誤りが無いことが確認されると、システムは自動的にユーザーをシステムのメイン画面へジャンプさせます。そうでなければ、システムはパスワード間違いやアカウントが存在しないなどのログイン失敗の具体的な原因をユーザーにフィードバックします。

图片生成及处理系统

登录

用户名

还没有账号？

去注册

密码

登录

4.2.2 メインページ

ユーザーのログイン成功後、メイン画面に入ります。左上には自分のアカウント名が表示され、ログアウトも可能です。左側にはナビゲーションバーがあり、各機能へ簡単にアクセスできます。右側には各種ボタンを通じて対応する機能に入ることができます。



4.2.3 文字抽出

の機能は画像をアップロードし、画像内の文字を自動認識して出力することができ、囲まれた文字の位置を表示します。同時に、メインページに戻るボタンも内蔵されています。

文字提取（OCR）

请选择图片用于 OCR (png/jpg/jpeg/bmp)

 Drag and drop file here
Limit 200MB per file • PNG, JPG, JPEG, BMP

Browse files

 屏幕截图 2025-09-10 160452.png 11.3KB

×

 原图

 带识别框的图片

上传并识别

识别结果（原文）

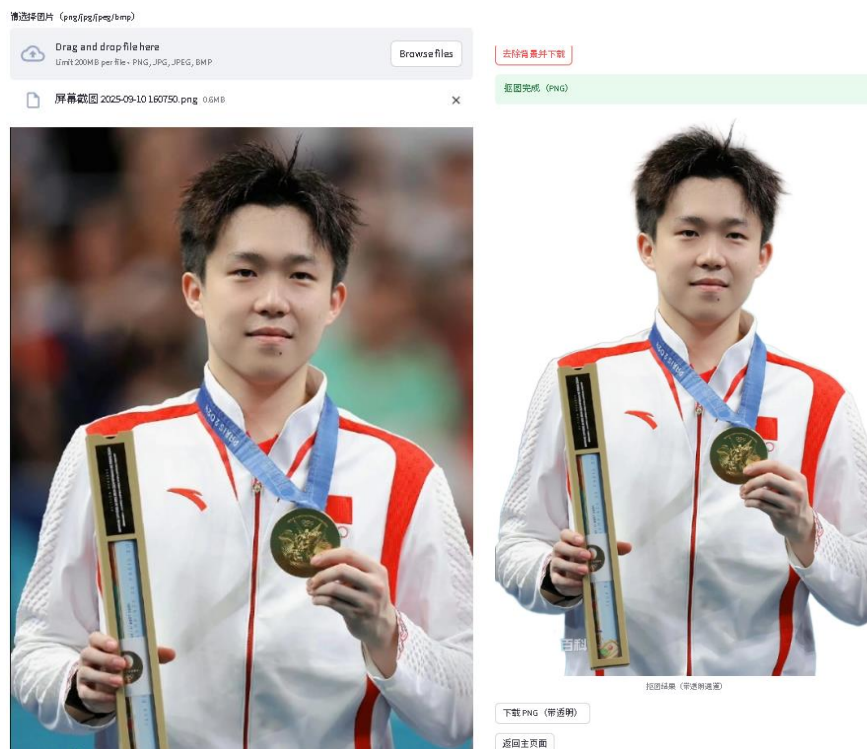
识别到的文本

该小说主要讲述了孙悟空出世并寻菩提祖师学艺及大闹天宫后。与猪八戒。沙僧和白龙马一同护送唐僧取经。路上历经险阻。降妖除魔。渡过了九九八十一难。成功到达大雷音寺。向如来佛祖求得《三藏真经》。最后五圣成真的故事。该小说以“玄奘取经”这一历史事件为蓝本；

返回主页面

4.2.4 画像切り抜き

ユーザーは人物を含む画像をアップロードでき、システムは自動的に背景を除去します。そして、人物のみのPNGをダウンロードすることが可能です。同時に、メインページに戻るボタンも内蔵されています。



4.2.5 テキスト画像生成

本機能では、ユーザーは自分の記述を入力し、画像サイズと数量を選択できます。さらに入力の拡張や透かし（ウォーターマーク）の追加などを選択でき、その後自分のAPIキーを入力すれば、画像を生成できます。同様にメインページへ戻るボタンも内蔵されています。

文生图 / 生成图片（使用 Dashscope / qwen-image）

填写 prompt，选择参数，生成图片。

Prompt（中文/英文均可）

一座山

生成数量 (n)

1

图片尺寸

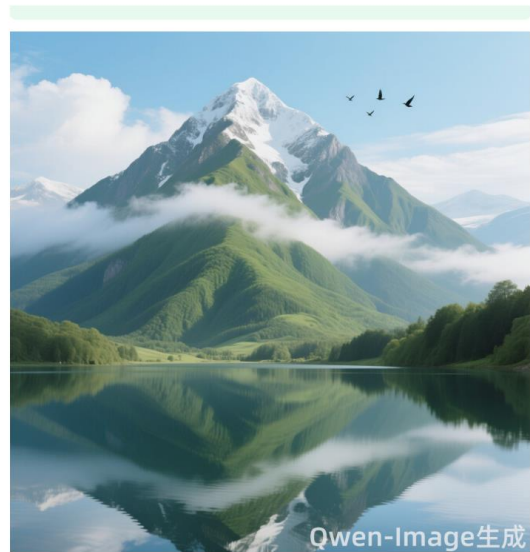
1328*1328

☒ 自动扩展 Prompt
(prompt_extend)

☒ 添加水印 (watermark)

填入你的api_key

生成图片（后端）



下载图片 (.png)

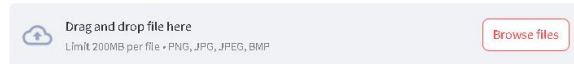
返回主页面

4.2.6 トリミングと解像度変更

これら2つの機能は、画像に対してトリミングや解像度の調整などを行うことができ、調整後のプレビュー画像を表示できます。同時にメインページへ戻るボタンもあります。

上传图片 -> 使用滑块选择裁剪区域（百分比）-> 预览 -> 裁剪并下载（后端裁剪并返回PNG）。

请选择图片（png/jpeg/bmp）

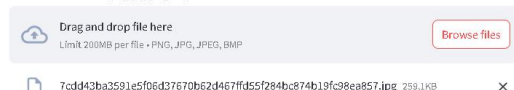


原图尺寸：1328 × 1328（像素）



返回主页面

请选择图片用于调整（png/jpeg/bmp）



原图尺寸：1328 × 1328（像素）



5 テキスト認識、人物認識、テキスト画像生成

本節では、EasyOCR、MODNet、Dashescope ImageSynthesisの動作原理について詳細に紹介します。

5.1 EasyOCR

EasyOCRはエンドツーエンドの光学文字認識ライブラリであり、そのワークフローは画像の前処理から始まります。ノイズ除去、グレースケール化、透視/幾何補正を通じてテキストの可読性を高め、その後テキスト領域の検出、領域ごとの認識、および後処理を行います。検出段階では通常、畳み込みニューラルネットワークに基づくテキスト検出器（例えばCRAFT類の手法）を採用して任意の方向とサイズのテキストブロックを特定し、必要に応じてテキスト行に対して薄板スプライン（TPS）やその他の変形補正を実施して文字配列を復元します。認識段階では、切り出された各テキスト行に対して深層認識モデルを実行します。まず畳み込みネットワークでシリアル化された特徴を抽出し、次にシーケンスモデル（双方向RNN/LSTMやアテンションベースのデコーダなど）を通じて時間次元でのモデリングを行い、学習時にはCTC損失やアテンションメカニズムを用いて整列とデコードを行い、推論時には貪欲法やビームサーチを組み合わせることで最終的なテキストを産出し、信頼度スコアを与えます。

多言語とエンジニアリング実践のレベルでは、EasyOCRはサブワードやグリフに基づく出力表現を採用して数十種類の文字セットをサポートし、事前学習済みの畳み込みバックボーンと言語固有の辞書/後処理戦略に依存して正確率を向上させています。フレームワーク全体はPyTorchなどのディープラーニングライブラリで実装され、リアルタイムまたは一括推論に向けて速度とメモリのバランスが調整されており、ドキュメント画像、ストリートビューテキスト、自然シーンの文字認識などのタスクにおいて良好な汎用性と拡張性を備えています。

5.2 MODNet

MODNet (Matting Objective Decomposition Network) は、人物切り抜き向けの軽量化された trimap-free リアルタイムモデルです。単一の RGB 画像のみを入力とし、エンドツーエンドの単一ネットワークを通じて高品質な alpha-matte を生成できます。その設計目標は、精細なエッジの詳細を保持しつつ、対話やリアルタイムプレビューの速度要求を満たすことです。

その核心理念は目的分解 (objective decomposition) です。全体的な切り抜きタスクを、意味推定 (Semantic estimation)、詳細予測 (Detail prediction)、意味-詳細融合 (Semantic-Detail fusion) の3つの相互に制約するサブ目標に分解し、対応する3つのブランチ (低解像度意味ブランチ S、高精細詳細ブランチ D、融合ブランチ F) を用いて単一ネットワーク内で並行して最適化します。これにより中間表現を共有し、誤差の連鎖を減少させます。効率と堅牢性を向上させるため、著者は効率的な Atrous Spatial Pyramid Pooling (e-ASPP) を導入して高速なマルチスケール意味融合を行い、サブ目標一貫性に基づく自己教師あり戦略 (SOC) を提案して実シーン下でのドメイン偏差問題を緩和しました。モデルはエンドツーエンドで学習可能で、512×512 入力下で約 67 FPS (1080Ti) に達し、公開ベンチマークにおいて同時代の trimap なし手法よりも著しく優れていました。コードと事前学習済みモデルはオープンソース化され、多くの推論フレームワーク (MobileNetV2 に基づく軽量化バリエーションなど) に移植されています。

全体として、MODNet は「分解—並行制約—高効率融合」というエンジニアリング設計を通じて、精度、効率、実用性の間で良好なバランスを実現しており、低遅延の人物切り抜きを必要とするリアルタイム応用シーン (ビデオ通話、カメラプレビュー、モバイルでの背景置換など) に適しています。

5.3 Dashscope ImageSynthesis

DashScope ImageSynthesis の「文生図」は、アリババクラウドの大規模モデルサービス (Model Studio / DashScope SDK) が外部に公開しているテキストから画像への生成能力を持つクラウドインターフェースです。ユーザーは REST API または

DashScope SDKを通じて自然言語のプロンプトおよび生成パラメータを送信し、プラットフォームがバックエンドのモデルと計算能力のスケジューリングを担当し、生成結果（通常はOSS画像リンクの形式）を返します。

生成フローにおいて、プラットフォームはまず入力プロンプトに対してテキストエンコーディングと意味表現を行い、その後この条件情報をバックエンドの画像生成モデル（現在のエコシステムで一般的な実装には、拡散モデルに基づく潜在ノイズ除去プロセスやdiffusion-transformer混合アーキテクチャなどが含まれる）に送ります。モデルは逆ノイズ除去/サンプリングの反復の中で段階的にピクセルまたは潜在表現を構築し、その後デコード、オプションの超解像/美的後処理、セキュリティ審査を経て、最終的にダウンロード可能な画像結果を産出します。

このサービスは開発者のために底層の学習とデプロイの複雑さを隠蔽し、複数のプリセット文生図モデル（通義千問/通義万相、Stable Diffusionシリーズなど）とパラメータ調整（解像度、スタイル、サンプリングステップ数、ランダムシードなど）をサポートし、呼び出し量、流量制限、結果保全の面でドキュメントと運用保守の制約を提供しており、エンジニアリング接続と製品化デプロイにより適しています。

5.4 大規模モデルによるテキスト画像生成

文生図（text-to-image）は、自然言語の記述を直接画像にマッピングする技術を目指します。現代の大規模モデルは通常、このタスクを「テキスト理解 → 条件付き生成 → デコードと後処理」の3段階に分割しており、核心はマルチモーダルアライメント（text-image alignment）と段階的生成（iterative denoising / autoregressive）メカニズムに依存しています。

テキスト前処理とエンコーディング

入力テキストに対してトークン化（**tokenization**）を行い、事前学習済みテキストエンコーダー（**Transformer**ベースのテキスト埋め込み器など、一般的なのは**CLIP-style**や自己学習のテキストエンコーダー）に送ります。出力は言語の意味に対応するベクトルシーケンスまたはグローバル意味ベクトルであり、生成モデルの条件情報となります。

条件情報の融合（**conditioning**）

生成ネットワーク（通常は**U-Net / Transformer**）内部で、クロスアテンション（**cross-attention**）または条件埋め込みを通じてテキストベクトルをモデルに注入し、生成プロセスの各ステップでテキストの意味を参照できるようにします。サンプリング時に**classifier-free guidance**などの手法を採用し、条件信号を増幅してテキストへの従順性を高めます。

生成の核心：ノイズモデリングとサンプリング

現代の主流な手法の多くは、拡散モデル（**diffusion models**）またはその潜在空間バリエーション（**latent diffusion**）です。学習時、モデルはノイズから画像を復元する方法（ノイズ予測/再構成損失の最小化）を学習します。サンプリング時は逆拡散（マルチステップノイズ除去）またはその高速化バリエーション（**DDIM**、**PLMS**など）を実行し、ランダムノイズを条件記述に合致する画像の潜在表現へと段階的に変換します。

デコードと超解像

潜在空間（**latent space**）で生成する場合、事前学習済みの**VAE/decoder**を通じて潜在表現をピクセル画像にデコードする必要があります。解像度と細部の一貫性を高めるために、超解像/詳細強化モジュール（**upsampler**）や後処理ネットワークを接続することがよくあります。

学習の詳細とマルチタスク目標

学習には大規模な画像-キャプションペア（**image-caption pairs**）を使用し、主な損失には拡散再構成損失、対照的アライメント損失（オプション）、およびテキスト

-画像一貫性正則化項が含まれます。データの選別、拡張、重複排除、ロングテールサンプルの処理は品質にとって極めて重要です。

安全性と制約

デプロイ段階でコンテンツ審査、セキュリティフィルタリング（暴力、ヘイト、個人の肖像など敏感なコンテンツの検出）、および著作権/スタイル制約戦略を加え、不適切な生成を防止します。

まとめ

文生図の技術スタックは、強力な意味エンコーディング + 条件付き生成ネットワーク + 段階的サンプリングと高品質デコードで構成されています。拡散モデルとクロスアテンションメカニズムは複雑な自然言語記述への忠実度を著しく向上させました。一方、データ品質、条件融合戦略（指導強度など）、およびバックエンドのデコードとセキュリティ戦略が、最終画像のリアリズムとコンプライアンスを決定します。将来の改善方向には、主に意味-視覚アライメント精度の向上、サンプリングコストの削減、およびより堅牢なスタイル/著作権管理とセキュリティ戦略が含まれます。

6 評価

本システムは、文生図（DashScope ImageSynthesis）、トライマップ不要人物切り抜き（MODNet）、多言語文字認識（EasyOCR）の3大能力を同一プラットフォーム上でエンジニアリング的に統合することに成功し、ユーザー認証、画像生成、前景分離から文字認識とエクスポートまでの閉ループワークフローを構築しました。システム設計は実用性と拡張性のバランスを体現しています。モジュール化されたインターフェース、テンプレート化されたプロンプト（文生図の安定性向上）、切り抜きとOCRの後処理およびキャッシュ戦略（遅延低減）、そして多形式エクスポートと透明背景サポートは、いずれもユーザー体験とデプロイ効率の向上に寄与しています。論文では重要モジュールの技術原理とエンジニアリング最適化についても明確な記述がなされており、システムが正確性、応答速度、使いやすさの面でエンジニアリング的に使用可能なレベルに達していることを示しています。

同時に、システムには注視すべき改善点も存在します。クラウドサービス（DashScope）への依存がもたらすコスト、遅延、プライバシーリスクに対して、代替/ローカルフォールバック戦略を策定すべきです。モジュール間の意味的一貫性（例えばOCR解析エラーが生成結果に及ぼす連鎖的影響）は、より多くの自動補正とヒューマン・マシン・インタラクションのメカニズムを通じて緩和する必要があります。切り抜きは極端な照明、複雑な背景、あるいは細かい髪の毛のディテールにおいて依然として欠陥がある可能性があり、より完全な定量的評価（IoU、エッジ誤差、OCR認識率）と大規模なユーザーユーザビリティテストを補足して汎化性を検証することを推奨します。全体として、本システムはアーキテクチャが明確で機能カバレッジも全面的であり、製品プロトタイプや業界ツールの技術基盤として適しています。厳格な性能ベンチマーク、プライバシー保護、フォールトトレランス（耐障害性）メカニズムを補足し、継続的な監視とモデルのイテレーションフローを導入すれば、より高い本番稼働への準備度と商用展開の潜在能力を備えることになるでしょう。

7 まとめと展望

本文では、DashScope ImageSynthesisのテキスト画像生成機能、MODNetのトライマップ不要人物切り抜き能力、EasyOCRの多言語文字認識モジュールを同一プラットフォームに統合した統合型画像処理システムを提案し実装しました。モジュール化設計と統一インターフェースを通じて、ユーザー認証、画像生成、前景分離から文字認識とエクスポートまでの閉ループプロセスを構築しました。システムはエンジニアリング実装レベルで複数の最適化を行いました。切り抜きとOCRのために安定した推論パイプラインと後処理メカニズムを構築し、テンプレート化されたプロンプトと指令解析を採用して文生図の一貫性を高め、キャッシュと軽量化デプロイを導入して遅延を削減し、透明背景、多形式エクスポート、設定可能なパラメータを提供して実用性と拡張性を強化しました。小規模テストでは、システムが応答速度、認識/切り抜き精度、ユーザーユーザビリティの面で良好なパフォーマンスを備えていることを示しており、製品プロトタイプや業界ソリューションとしての基礎条件をすでに備えています。

将来の展望として、いくつかの重要な方向で深化と改善を図る必要があります。まず、クラウドサービスへの依存に対し、ローカルまたはハイブリッド推論のフォールバック戦略を導入し、コスト削減、遅延緩和、データプライバシー保護の向上を図る必要があります。次に、モジュール間の意味的一貫性検証と自動修正フローを強化し、OCR結果と切り抜きマスクを生成後処理のフィードバック信号として利用し、連鎖誤差を減らして最終出力の意味的正確度を向上させます。さらに、体系的な定量的評価（IoU、エッジ誤差、OCR再現率/適合率、エンドツーエンド遅延など）を補足し、大規模なユーザーユーザビリティ実験を展開して、汎化性と製品化の実現可能性を検証することを推奨します。技術ルートの面では、Transformerなどの現代的構造を導入して高解像度切り抜きと細部復元能力を向上させることや、差分プライバシーや連合学習などのメカニズムを研究してプライバシー保護を強化することを探索できます。全体として、本システムは文生図、切り抜き、OCRの協調応用に対して明確なエンジニアリングソリューションと初期検証を提供しており、

厳格なベンチマークテスト、プライバシーおよびフォールトトレランスメカニズム、そして継続的なイテレーションフローと組み合わせることで、本番レベル、商用化の方向へ着実に推進できると期待されます。