

基于stable diffusion的数据增强及数据标注

在计算机视觉领域，图像生成指根据模型或算法生成新的图像。将图像生成技术应用
于实际产业能有效节省生产时间和人工成本，提高工作效率。目前该技术已被应
用于艺术品生成、广告设计、游戏开发、机器人视觉和虚拟现实等领域，具有广阔
的应用前景。

1 扩散模型与Stable Diffusion

1.1 扩散模型

我们根据使用技术的不同，可将图像生成方法分为传统方法和基于深度学习的方法。传统方法包括基于纹理合成和基于马尔可夫随机场等方法。然而，传统方法的生成效果欠佳，且通常仅适用于特定的任务场景。基于深度学习的方法根据学习类型可分为生成对抗网络、自回归模型、变分自编码器、流模型、能量模型和扩散模型等。与传统方法相比，基于深度学习的图像生成方法生成的图像质量更好，能满足实际应用需求。在实际应用中大多对生成内容有一定要求，条件引导的图像生成方法能增强对生成过程的控制，使生成内容向目标要求迈进；与其他生成模型相比，扩散模型在条件引导的图像生成中具有生成质量高和多样性强等优点，亦展现出巨大的发展潜力。

1.1.1 扩散模型的定义

扩散模型的目标可以总结为逆转数据逐渐退化的过程，包括符合马尔科夫链的正向过程和逆扩散过程，如图 1 所示。正向过程，在原始数据中逐步添加噪声，使其逐渐退化为几乎各向同性的高斯噪声，破坏原始数据；逆扩散过程，通过神经网络学习从高斯噪声中恢复原始数据。需要指出的是，逆扩散过程输入和输出的数据维度需保持不变，由于 U-Net 符合此要求且开销较小，被广泛用于构建扩散模型的去噪网络。

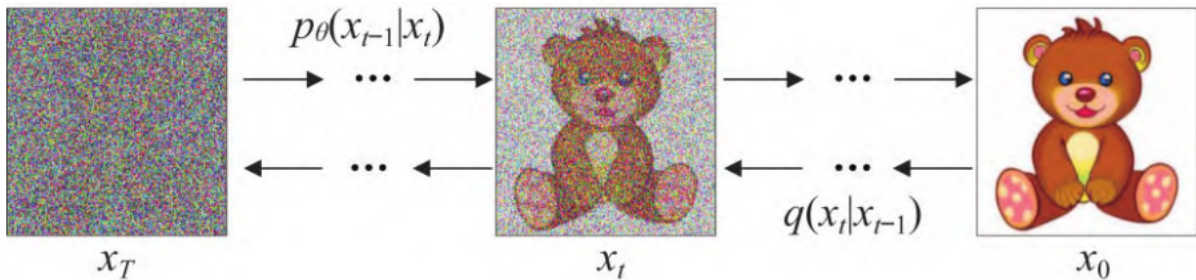


图 1 扩散模型的正向和逆扩散过程

根据定义方法，扩散模型可分为去噪扩散概率模型（denoising diffusion probabilistic models，

DDPMs）、分数生成模型（score-based generative models，SGMs）和随机微分方程（stochastic differential equations，SDEs）3类。

1.1.2 分类器引导

在基于扩散模型的图像生成中，显式分类器引导、隐式分类器引导和多模态引导均能提高图像生成的灵活性和质量。

1 显式分类器

虽然 DDPM和改进的 DDPM等方法均能生成逼真的图像，但在 FID 等指标上与基于生成对抗网络的主流方法仍存在一定差距。为了解决此问题，DHARIWAL 等引入了显式分类器引导（classifier guidance），通过在图像生成过程中注入类别标签信息增强对模型的约束，从而提高图像生成质量。基于该技术，由扩散模型生成的图像 FID 等指标超过了基于生成对抗网络的 BigGAN。

通过去噪网络构建 $p(x_{t-1}|x_t)$ 是扩散模型的关键，利用贝叶斯公式显式分类器将预训练无条件图像生成模型中的 $p(x_{t-1}|x_t)$ 转化为 $p(x_{t-1}|x_t, y)$ ：

$$p(x_{t-1}|x_t, y) = \frac{p(x_{t-1}|x_t)p(y|x_{t-1}, x_t)}{p(y|x_t)}$$

其中 y 表示给定的条件。由于该过程仍基于预训练无条件图像生成模型的生成过程 $p(x_{t-1}|x_t)$ ，所以不需要训练新模型。

经计算， $p(x_{t-1}|x_t, y)$ 近似于 $\mathcal{N}(x_{t-1}; \mu(x_t) + \sigma_t^2 \nabla_{x_t} \lg p(y|x_t), \sigma_t^2 I)$ ，其中中间项就是无条件生成过程的新增项。

虽然显式分类器的成本较低，但存在以下缺陷：（1）需额外训练一个分类器，增加了模型的复杂性；（2）分类器在一定程度上决定了模型的上限，优化该分类器也具有一定的难度；（3）显式分类器引导会破坏生成结果的多样性等。

2 隐式分类器

受显式分类器的启发，HO 等提出了隐式分类器引导（classifier-free guidance）。与显式分类器不同，隐式分类器可直接将转换函数 $p(x_{t-1}|x_t, y)$ 定义为， $\mathcal{N}(x_{t-1}; \mu(x, y), \sigma_t^2 I)$ 这是两者的本质差别。

隐式分类器无需训练额外的分类器，而是训练条件生成和无条件生成 2 类模型。在训练模型时，引入额外的条件 y ，在不需要条件时用空值代替条件 y ，从而降低训练难度，但大大增加了训练成本。目前的图像生成大模型 **Glide**、**Stable Diffusion**、**DALLE2** 和 **Imagen** 等均采用隐式分类器，生成效果惊人。

3 多模态引导

基于 **CLIP** 的多模态引导被应用于条件引导的图像生成、文本引导的图像编辑和文本引导的图像转换等任务。预训练的 **CLIP** 模型能很好地衡量图像与文本之间的相似程度特性，有助于模型构建

$$\mathcal{L}(x_{\text{gen}}, y_{\text{tar}}; x_{\text{ref}}, y_{\text{ref}}) = 1 - \frac{\langle \Delta I, \Delta T \rangle}{\|\Delta I\| \|\Delta T\|}$$

其中 $\Delta T = E_T(y_{\text{tar}}) - E_T(y_{\text{ref}})$ 为文本之间的向量差值, $\Delta I = E_I(y_{\text{tar}}) - E_I(y_{\text{ref}})$ 为图像之间的向量差

值, E_T 和 E_I 分别为 **CLIP** 模型的文本编码器和图像编码器。

由于 **CLIP** 模型在生成过程中计算每步损失均极其耗时，随着隐式分类器的提出，基于 **CLIP** 的多模态引导逐渐被淘汰。

综上，条件引导图像生成的条件和对应方法如图 2 所示。

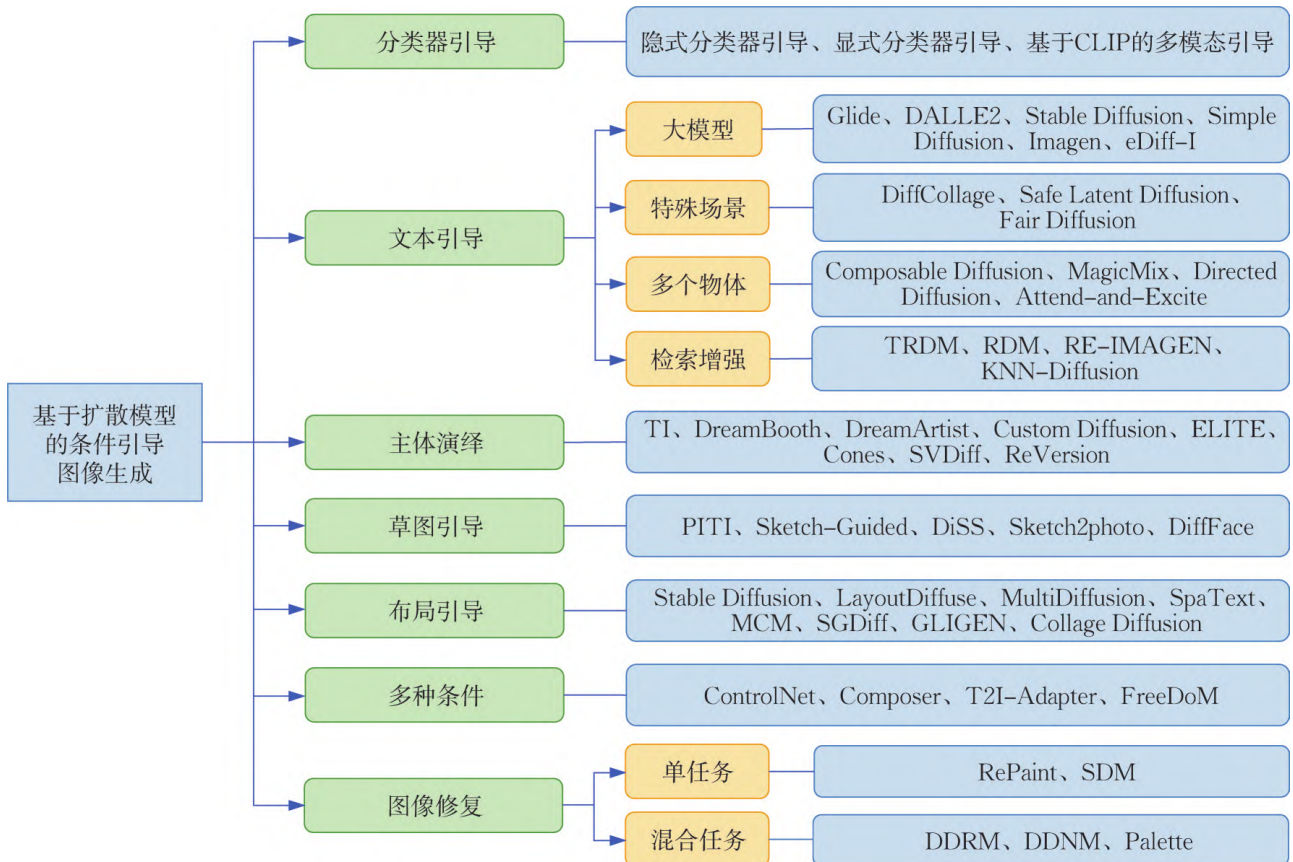


图 2 条件引导图像生成的条件和对应方法

1.1.2 条件引导图像生成的应用

随着扩散模型的发展，近几年出现了许多条件引导的图像生成模型，其中与文本相关的条件引导

生成模型占多数，这与文本数据易获取、数量多有关，且利用文本控制图像生成的过程相对简单。此外，还有针对特定任务的生成模型，如基于草图、布局的生成模型等。

1 基于文本引导的图像生成模型

计算机视觉和自然语言处理是深度学习的两个热门方向，为此有大量的图像数据和文本数据被保存，其中有多组为开源的大规模图像文本数据集。与无条件图像生成方法相比，基于文本引导的图像生成方法能显著提高图像生成的可控制性，不仅自动化程度高、灵活性强，还能在节省时间成本的同时应对复杂的场景。下面主要介绍大型文本生成图像模型、特定场景下的文本生成图像模型、生成包含多个物体图像的文本生成图像模型和基于检索增强的文本生成图像模型。

目前实际应用最多的文本生成图像模型主要有Stable Diffusion、Midjourney 和 DALL·E2。BORJI A.在真实场景下对上述模型生成人脸图像的能力进行了对比。结果表明，Midjourney生成的图像具有明显的动漫和艺术风格，真实性最差，Stable Diffusion 生成的图像真实性优于 DALL·E2。FID 指标，Stable Diffusion 最好，其次是 DALL·E2，最后是Midjourney。

2 对图像的主体内容进行演绎的图像生成模型

尽管现有的大型文本生成图像模型生成的图像具有良好的真实性和多样性，但无法在给定参考图像的情况下对参考图像中的主体进行相关的演绎。Textual Inversion在不改变图像主体基本属性的情况下，根据文本对图像中的主体进行了创造性演绎，首先通过隐向量空间的文本编码器学习新的概念，然后根据文本所包含的特定概念实现对图像的精细控制。

3 以草图为条件的图像生成模型

基于扩散模型的条件图像生成方法能够产生具有显著多样性和真实感的图像，然而，大多方法只允许对标签或文本提示进行调节，限制了对最终结果的控制强度。因此，出现了一些以草图为条件的图像生成方法。

1.2 Stable diffusion

扩散模型最大的问题是它的时间成本和经济成本都极其“昂贵”。Stable Diffusion的出现就是为了解决上述问题。如果我们想要生成一张 $1024 \times 1024 \times 3$ 尺寸的图像，U-Net 会使用 $1024 \times 1024 \times 3$ 尺寸的噪声，然后从中生成图像。这里做一步扩散的计算量就很大，更别说要循环迭代多次直到100%。一个解决

方法是將大圖片拆分為若干小分辨率的圖片進行訓練，然後再使用一個額外的神經網絡來產生更大分辨率的圖像（超分辨率擴散）。

2021年發布的Latent Diffusion模型給出了不一樣的方法。Latent Diffusion模型不直接在操作圖像，而是在潛在空間中進行操作。通過將原始數據編碼到更小的空間中，讓U-Net可以在低維表示上添加和刪除噪聲。

1.2.1 潛在空間(Latent Space)

潛在空間簡單的說是对壓縮數據的表示。所謂壓縮指的是用比原始表示更小的數位來編碼信息的过程。比如我們用一個顏色通道（黑白灰）來表示原來由RGB三原色構成的圖片，此時每個像素點的顏色向量由3維變成了1維度。維度降低會丟失一部分信息，然而在某些情況下，降維不是件壞事。通過降維我們可以過濾掉一些不太重要的信息，只保留最重要的信息。假設我們像通過全連接的卷積神經網絡訓練一個圖像分類模型。當我們說模型在學習時，我們的意思是它在學習神經網絡每一層的特定屬性，比如邊緣、角度、形狀等.....每當模型使用數據（已經存在的圖像）學習時，都會將圖像的尺寸先減小再恢復到原始尺寸。最後，模型使用解碼器從壓縮數據中重建圖像，同時學習之前的所有相關信息。因此，空間變小，以便提取和保留最重要的屬性。這就是潛在空間適用於擴散模型的原因。

1.2.2 Latent Diffusion

“潛在擴散模型”（Latent Diffusion Model）將GAN的感知能力、擴散模型的細節保存能力和Transformer的語義能力三者結合，創造出比上述所有模型更穩健和高效的生成模型。與其他方法相比，Latent Diffusion不僅節省了內存，而且生成的圖像保持了多樣性和高細節度，同時圖像還保留了數據的語義結構。

語義壓縮

在學習的第二階段，圖像生成方法必須能夠捕獲數據中存在的語義結構。這種概念和語義結構提供了圖像中各種對象的上下文和相互關係的保存。Transformer擅長捕捉文本和圖像中的語義結構。Transformer的泛化能力和擴散模型的細節保存能力相結合，提供了兩全其美的方法，並提供了一種生成細粒度的高度細節圖像的方法，同時保留圖像中的語義結構。

1.2.3 感知損失

潛在擴散模型中的自動編碼器通過將數據投影到潛在空間來捕獲數據的感知結構。論文作者使用一種特殊的損失函數來訓練這種稱為“感知損失”的自動編碼器。該損失函數確保重建限制在圖像流形內，並減少使用像素空間損失（例如 L1/L2 損失）時出現的模糊。

1.2.4 扩散损失

扩散模型通过从正态分布变量中逐步去除噪声来学习数据分布。换句话说，扩散模型使用长度为 T 的反向马尔可夫链。这也意味着扩散模型可以建模为时间步长为 $t = 1, \dots, T$ 的一系列“ T ”去噪自动编码器。由下方公式中的 ϵ_θ 表示：

$$L_{DM} = E_{x, \epsilon \sim N(0,1), t} [\| \epsilon - \epsilon_\theta(x_t, t) \|_2^2]$$

公式(1)给出了扩散模型的损失函数。在潜在扩散模型中，损失函数取决于潜在向量而不是像素空间。我们将像素空间元素 x 替换成潜在向量 $\epsilon(x)$ ，将 t 时间的状态 x_t 替换为去噪 U-Net 在时间 t 的潜在状态 z_t ，即可得到潜在扩散模型的损失函数，见公式(2)：

$$L_{LDM} := E_{\epsilon(x), \epsilon \sim N(0,1), t} [\| \epsilon - \epsilon_\theta(z_t, t), \tau_\theta(y) \|_2^2]$$

将公式(2)写成条件损失函数，得到公式(3)：

$$L_{LDM} := E_{\epsilon(x), y, \epsilon \sim N(0,1), t} [\| \epsilon - \epsilon_\theta(z_t, t), \tau_\theta(y) \|_2^2]$$

1.2.5 条件扩散

扩散模型是依赖于先验的条件模型。在图像生成任务中，先验通常是文本、图像或语义图。为了获得先验的潜在表示，需要使用转换器（例如 CLIP）将文本/图像嵌入到潜在向量 τ 中。因此，最终的损失函数不仅取决于原始图像的潜在空间，还取决于条件的潜在嵌入。

1.2.6 注意力机制**

潜在扩散模型的主干是具有稀疏连接的 U-Net 自动编码器，提供交叉注意力机制²。

Transformer 网络将条件文本/图像编码为潜在嵌入，后者又通过交叉注意力层映射到 U-Net 的中间层。这个交叉注意力层实现了注意力 $(Q, K, V) = \text{softmax}(QKT/\sqrt{d})V$ 其中 Q 、 K 和 V 是可学习的投影矩阵

1.2.7 文本-图像合成

在 Python 实现中，我们可以使用使用 LDM v4 的最新官方实现来生成图像。在文本到图像的合成中，潜在扩散模型使用预训练的 CLIP 模型³，该模型为文本和图像等多种模态提供基于 Transformer 的通用嵌入。然后将 Transformer 模型的输出输入到称为“diffusers”的潜在扩散模型 Python API，同时还可以设置一些参数（例如，扩散步数、随机数种子、图像大小等）。

1.2.8 图像-图像合成

相同的方法同样适用于图像到图像的合成，不同的是需要输入样本图像作为参考图像。生成的图像在语义和视觉上与作为参考给出的图像相似。这个过程在概念上类似于基于样式的 GAN 模型，但它在保留图像的语义结构方面做得更好。

1.2.9 整体架构

上面介绍了潜在扩散模型各个主要技术部分，下面我们将它们合成一个整体，看一下潜在扩散模型的完整工作流程。

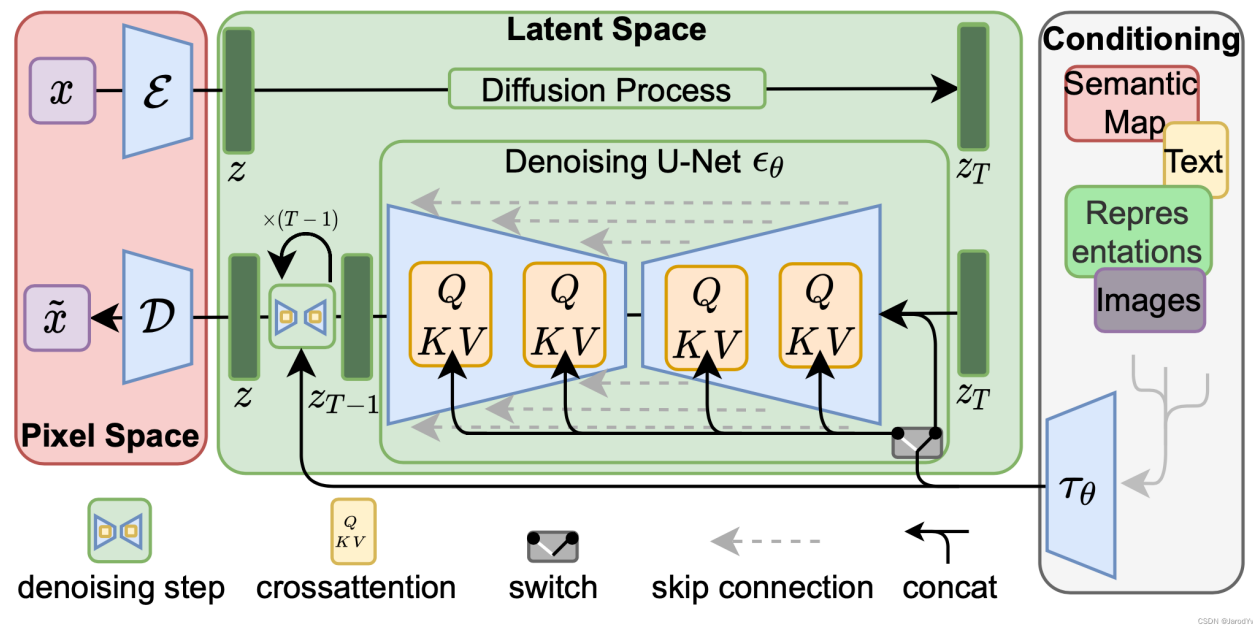


图 3 Stable Diffusion整体架构

上图中 x 表示输入图像， \tilde{x} 表示生成的图像； \mathcal{E} 是编码器， \mathcal{D} 是解码器，二者共同构成了感知压缩； z 是潜在向量； z_T 是增加噪声后的潜在向量； τ_θ 是文本/图像的编码器（比如 Transformer 或 CLIP），实现了语义压缩。

2 实验过程

2.1 Stable Diffusion 的本地部署

本地部署 Stable Diffusion 最简单的方法是使用 Stable Diffusion Web Ui。

Stable Diffusion Web Ui 是一套无代码、可视化的 Stable Diffusion 集成运行环境。它将 Stable Diffusion 的安装部署集成打包，提供一键安装脚本，并提供 Web 界面操作界面，极大简化了 Stable Diffusion 的操作和使用，让没有不懂代码的小白也能轻松上手使用 Stable Diffusion 模型。

我们采用整合包直接安装Stable Diffusion Web Ui（在这里感谢秋葉aaaki），安装后只需下载启动依赖即可。

2.2 使用Stable Diffusion进行数据增强

2.2.1 提示词

prompt: 主要是对于图像进行描述。**prompt**对Stable Diffusion图像生成质量至关重要，因此如果想生成高质量图片，一定要在提示设计上下功夫。一个好的提示需要详细和具体。

Negative prompt: 主要是告诉模型我不想要什么样的风格或元素；

prompt语法

为了产生具有特定风格的图像，必须以特定格式提供文本提示。这通常需要添加提示修饰符或添加更多关键字或关键短语来实现。下面为大家介绍一下Stable Diffusion的**prompt**语法规则。

Stable Diffusion提示文本中的关键字或关键短语通过半角逗号分割，一般越靠前权重越高。我们可以通过提示修饰符来认为修改权重。

- (tag): 增加权重5%
- [tag]: 降低权重5%
- (tag: weight): 设置具体权重值

括号可以嵌套使用，例如：(tag)的权重为 $1 \times 1.05 = 1.05$ ((tag))的权重为 $1 \times 1.05 \times 1.05 = 1.10255$ 。同理[tag]的权重为 $\frac{1}{1.05} = 0.952$ [[tag]]的权重为 $\frac{1}{1.05^2} = 0.907$

- [tag1 | tag2]: 将tag1和tag2混合；
- {tag1 | tag2 | tag3}: 从标签集合中随机选择一个标签；
- [tag1 : tag2 : 0.5]: 表示先用tag1生成，当生成进程到50%时，改用tag2生成；如果输入整数的话表示步长，比如10，意思是生成10步后改用tag2；
- <lora:filename:multiplier>: LoRA模型引用语法

2.2.2 Stable Diffusion 模型

与DALL-E和Midjourney相比，Stable Diffusion最大的优势是开源，这就意味着Stable Diffusion的潜力巨大、发展飞快。Stable Diffusion已经跟很多工具和平台进行了集成，且可用预训练模型数量众多（参见Stable Diffusion资源列表）。正是由于社区的活跃，使得Stable Diffusion在各种风格的图像生成上都有着出色的表现。

我们在civitai.com上选择合适的stable diffusion模型下载并导入，由于本次需要生成画风写实的图片，我们挑选了PicX_real作为Stable Diffusion模型。

2.3.3 调参

由于stable参数量众多，碍于篇幅，这里仅对调参过程中用到的参数进行介绍

迭代步数

模型生成图片的迭代步数，每多一次迭代都会给 AI 更多的机会去对比 prompt 和 当前结果，从而进一步调整图片。更高的步数需要花费更多的计算时间，但却不一定意味着会有更好的结果。当然迭代步数不足肯定会降低输出的图像质量。

采样方法

扩散去噪算法的采样模式，不同采样模式会带来不一样的效果，具体需要在实际使用中测试。

重绘程度

决定算法对图像内容的保留程度。0什么都不会改变，1会得到一个完全不同的图像；

Controlnet

调节完成以上内容后，只需点击批量生成即可开始生成图片

2.3 数据标注

使用基于百度飞桨的EasyDL平台对新扩充的数据进行标注。标注的类别主要分三类：分别是未满溢的垃圾桶、满溢的垃圾桶和垃圾（和论文中相同）。