

OTA: Optimal Transport Assignment for Object Detection

Abstract

Recent advances in label assignment in object detection mainly seek to independently define positive/negative training samples for each ground-truth (gt) object. In this paper, we innovatively revisit the label assignment from a global perspective and propose to formulate the assigning procedure as an Optimal Transport (OT) problem – a well-studied topic in Optimization Theory. Concretely, we define the unit transportation cost between each demander (anchor) and supplier (gt) pair as the weighted summation of their classification and regression losses. After formulation, finding the best assignment solution is converted to solve the optimal transport plan at minimal transportation costs, which can be solved via Sinkhorn-Knopp Iteration. On COCO, a single FCOS-ResNet-50 detector equipped with Optimal Transport Assignment (OTA) can reach 40.7% mAP under 1× scheduler, outperforming all other existing assigning methods. Extensive experiments conducted on COCO and CrowdHuman further validate the effectiveness of our proposed OTA, especially its superiority in crowd scenarios. The code is available at <https://github.com/Megvii-BaseDetection/OTA>.

目标检测中标签匹配的最新进展主要是寻求为每个真实框（gt）目标独立定义正/负的训练样本。在本文中，我们创新性地从全局角度重新审视了标签匹配，并建议将匹配过程表述为最佳传输（OT）问题--这是优化理论中一个被充分研究的课题。具体来说，**我们将每对需求者（anchor）和供应商（gt）之间的单位运输成本定义为其分类和回归损失的加权总和**。在制定之后，**寻找最佳分配方案被转换为解决运输成本最小的最佳运输计划**，这可以通过**Sinkhorn-Knopp迭代法**解决。在COCO上，一个配备了最优运输分配（OTA）的FCOS-ResNet-50检测器在1×调度器下可以达到40.7%的mAP，优于所有其他现有的分配方法。在COCO和CrowdHuman上进行的大量实验进一步验证了我们提出的OTA的有效性，尤其是在人群场景中的优越性。代码可在<https://github.com/Megvii-BaseDetection/OTA>。

1. Introduction

Current CNN-based object detectors [27, 30, 21, 47, 33, 29, 36] perform a dense prediction manner by predicting the classification (cls) labels and regression (reg) offsets for a set of pre-defined anchors¹. To train the detector, defining cls and reg targets for each anchor is a necessary procedure, which is called label assignment in object detection.

目前基于CNN的物体检测器[27, 30, 21, 47, 33, 29, 36]通过预测一组预先定义的锚点的分类（cls）标签和回归（reg）偏移量来执行密集的预测方式¹。为了训练检测器，为每个锚点定义cls和reg目标是一个必要的程序，这在物体检测中称为标签匹配。

Classical label assigning strategies commonly adopt pre-defined rules to match the ground-truth (gt) object or background for each anchor. For example, RetinaNet [21] adopts Intersection-over-Union (IoU) as its thresholding criterion for pos/neg anchors division. Anchor-free detectors like FCOS [38] treat the anchors within the center/bbox region of any gt object as the corresponding positives. Such static strategies ignore a fact that for objects with different sizes, shapes or occlusion condition, the appropriate positive/negative (pos/neg) division boundaries may vary.

经典的标签分配策略通常采用预先定义的规则来匹配每个锚点的地面真实（gt）对象或背景。例如，RetinaNet[21]采用交叉联合（IoU）作为其正/负锚点划分的阈值标准。像FCOS[38]这样的无锚检测器将任何gt对象的中心/方框区域内的锚视为相应的阳性。**这种静态策略忽略了一个事实，即对于具有不同尺寸、形状或遮挡条件的物体，适当的正/负（pos/neg）划分边界可能会有所不同。**

Motivated by this, many dynamic assignment strategies have been proposed. **ATSS** [47] proposes to **set the division boundary for each gt based on statistical characteristics**. Other recent advances [48, 19, 51, 16] **suggest that the predicted confidence scores of each anchor could be a proper indicator to design dynamic assigning strategies**, i.e., **high confidence anchors can be easily learned by the networks and thus be assigned to the related gt, while anchors with uncertain predictions should be considered as negatives**. Those strategies enable the detector to dynamically choose positive anchors for each individual gt object and achieve state-of-the-art performance.

受此启发，人们提出了许多动态分配策略。**ATSS**[47]提出，**根据统计学特征为每个gt设置划分边界**。其他最新进展[48, 19, 51, 16]提出，**每个锚的预测置信度分数可以作为设计动态分配策略的适当指标，即高置信度的锚可以很容易被网络学习，从而被分配到相关的gt，而预测不确定的锚应该被认为是负面的**。这些策略使检测器能够为每个单独的gt对象动态地选择正面的锚，并获得最先进的性能。

However, independently assigning pos/neg samples for each gt without context could be sub-optimal, just like the lack of context may lead to improper prediction. When dealing with ambiguous anchors (i.e., anchors that are qualified as positive samples for multiple gts simultaneously as seen in Fig. 1.), existing assignment strategies are heavily based on hand-crafted rules (e.g., Min Area [38], Max IoU [16, 21, 47]). We argue that assigning ambiguous anchors to any gt (or background) may introduce harmful gradients w.r.t. other gts. Hence the assignment for ambiguous anchors is non-trivial and requires further information beyond the local view. Thus a better assigning strategy should get rid of the convention of pursuing optimal assignment for each gt independently and turn to the ideology of global optimum, in other words, finding the global high confidence assignment for all gts in an image.

然而，在没有上下文的情况下，独立地为每个gt分配正/负样本可能是次优的，就像缺乏上下文可能导致不正确的预测一样。当处理模棱两可的锚（即图1中看到的同时被限定为多个gt的正样本的锚）时，现有的分配策略在很大程度上是基于手工制作的规则（如Min Area [38], Max IoU [16, 21, 47]）。我们认为，将模棱两可的锚点分配给任何gt（或背景），可能会对其他gt引入有害的梯度。因此，对模糊的锚点的分配是不难的，需要在本地视图之外的进一步信息。因此，一个更好的分配策略应该摆脱独立追求每个目标的最优分配的惯例，而转向全局最优的思想，换句话说，为图像中的所有目标找到全局的高置信度分配。

DeTR [3] is the first work that attempts to consider label assignment from global view. It replaces the detection head with transformer layers [39] and considers one-to-one assignment using the Hungarian algorithm that matches only one query for each gt with global minimum loss. However, for the CNN based detectors, as the networks often produce correlated scores to the neighboring regions around the object, each gt is assigned to many anchors (i.e., one to-many), which also benefits to training efficiency. In this one-to-many manner, it remains intact to assign labels with a global view.

DeTR[3]是第一个试图从全局角度考虑标签分配的工作。它用转换层[39]取代了检测头，并考虑使用**匈牙利算法**进行一对一的分配，**该算法对每个gt只匹配一个查询，且损失全局最小**。然而，对于基于CNN的检测器来说，由于网络经常对物体周围的邻近区域产生相关的分数，每个gt被分配给许多锚（即一对多），这也有利于训练效率。在这种一对多的方式中，它仍然是以全局的观点来分配标签的。

To achieve the global optimal assigning result under the one-to-many situation, we propose to formulate label assignment as an Optimal Transport (OT) problem – a special form of Linear Programming (LP) in Optimization Theory. Specifically, we define each gt as a supplier who supplies a certain number of labels, and define each anchor as a demander who needs one unit label. If an anchor receives sufficient amount of positive label from a certain gt, this anchor becomes one positive anchor for that gt. In this context, the number of positive labels each gt supplies can be interpreted as “how many positive anchors that gt needs for better convergence during the training process”. The unit transportation cost between each anchor-gt pair is defined as the weighted summation of their pair-wise cls and reg losses. Furthermore, as being negative

should also be considered for each anchor, we introduce another supplier – background who supplies negative labels to make up the rest of labels in need. The cost between background and a certain anchor is defined as their pair-wise classification loss only. After formulation, finding the best assignment solution is converted to solve the optimal transport plan, which can be quickly and efficiently solved by the off-the-shelf Sinkhorn-Knopp Iteration [5]. We name such an assigning strategy as Optimal Transport Assignment (OTA).

为了实现一对多情况下的全局最优分配结果，我们建议将标签分配表述为最优运输（OT）问题--优化理论中线性编程（LP）的一种特殊形式。具体来说，我们将每个gt定义为提供一定数量标签的供应商，并将每个锚定义为需要一个单位标签的需求者。如果一个锚从某个gt收到足够数量的正标签，这个锚就成为该gt的一个正锚。在这种情况下，每个gt提供的正向标签的数量可以解释为“在训练过程中，gt需要多少个正向锚来进行更好的收敛”。每个锚-gt对之间的单位运输成本被定义为它们成对的cls和reg损失的加权和。此外，由于每个锚的负值也应该被考虑在内，我们引入了另一个供应商--背景，他提供负值标签以弥补其余需要的标签。背景和某个锚之间的成本仅被定义为他们的配对分类损失。在制定之后，寻找最佳分配方案被转换为解决最优运输计划，这可以通过现成的Sinkhorn-Knopp迭代法[5]快速有效地解决。我们将这样的分配策略命名为最优运输分配（OTA）。

Comprehensive experiments are carried out on MS COCO [22] benchmark, and significant improvements from OTA demonstrate its advantage. OTA also achieves the SOTA performance among one-stage detectors on a crowded pedestrian detection dataset named CrowdHuman [35], showing OTA's generalization ability on different detection benchmarks.

在MS COCO[22]基准上进行了综合实验，OTA的明显改进显示了其优势。OTA还在名为CrowdHuman[35]的拥挤行人检测数据集上取得了单阶段检测器中的SOTA性能，显示了OTA在不同检测基准上的通用能力。

2. Related Work

2.1. Fixed Label Assignment

Determining which gt (or background) should each anchor been assigned to is a necessary procedure before training object detectors. Anchor-based detectors usually adopt IoU at a certain threshold as the assigning criterion. For example, RPN in Faster R-CNN [33] uses 0.7 and 0.3 as the positive and negative thresholds, respectively. When training the R-CNN module, the IoU threshold for pos/neg division is changed to 0.5. IoU based label assignment is proved effective and soon been adopted by many Faster R-CNN's variants like [2, 12, 20, 42, 26, 49, 37], as well as many one-stage detectors like [31, 32, 25, 27, 23, 21]. Recently, anchor-free detectors have drawn much attention because of their concision and high computational efficiency. Without anchor box, FCOS [38], Foveabox [17] and their precursors [30, 14, 46] directly assign anchor points around the center of objects as positive samples, showing promising detection performance. Another stream of anchor-free detectors [18, 8, 50, 45, 4] view each object as a single or a set of key-points. They share distinct characteristics from other detectors, hence will not be further discussed in our paper. Although detectors mentioned above are different in many aspects, as for label assignment, they all adopt a single fixed assigning criterion (e.g., a fixed region of the center area or IoU threshold) for objects of various sizes, shapes, and categories, etc, which may lead to sub-optimal assigning results.

在训练物体检测器之前，确定每个锚应该被分配到哪个gt（或背景）是一个必要的程序。基于锚点的检测器通常采用一定阈值的IoU作为分配标准。例如，Faster R-CNN[33]中的RPN分别使用0.7和0.3作为正负阈值。在训练R-CNN模块时，正/负划分的IoU阈值被改为0.5。基于IoU的标签分配被证明是有效的，并很快被许多Faster R-CNN的变体采用，如[2, 12, 20, 42, 26, 49, 37]，以及许多单阶段检测器如[31, 32, 25, 27, 23, 21]。最近，无锚检测器因其简洁和高计算效率而备受关注。在没有锚箱的情况下，FCOS[38]、Foveabox[17]及其前身[30、14、46]直接将物体中心周围的锚点作为阳性样本，显示出良好的检测性能。另一股无锚检测器[18, 8, 50, 45, 4]将每个物体视为单个或一组关键点。它们与其他检测

器有不同的特点，因此在我们的论文中不会进一步讨论。尽管上述检测器在很多方面都有所不同，但就标签分配而言，它们都对不同大小、形状和类别的物体采用了单一的固定分配标准（例如，中心区域的固定区域或IoU阈值），这可能导致次优的分配结果。

2.2. Dynamic Label Assignment

Many recent works try to make the label assigning procedure more adaptive, aiming to further improve the detection performance. Instead of using pre-defined anchors, GuidedAnchoring [40] generates anchors based on an anchor-free mechanism to better fit the distribution of various objects. MetaAnchor [44] proposes an anchor generation function to learn dynamic anchors from the arbitrary customized prior boxes. NoisyAnchors [19] proposes soft label and anchor re-weighting mechanisms based on classification and localization losses. FreeAnchor [48] constructs top-k anchor candidates for each gt based on IoU and then proposes a detection-customized likelihood to perform pos/neg division within each candidate set. ATSS [47] proposes an adaptive sample selection strategy that adopts mean+std of IoU values from a set of closest anchors for each gt as a pos/neg threshold. PAA [16] assumes that the distribution of joint loss for positive and negative samples follows the Gaussian distribution. Hence it uses GMM to fit the distribution of positive and negative samples, and then use the center of positive sample distribution as the pos/neg division boundary. AutoAssign [51] tackles label assignment in a fully data-driven manner by automatically determine the positives/negatives in both spatial and scale dimensions. These methods explore the optimal assigning strategy for individual objects, while failing to consider context information from a global perspective. DeTR [3] examines the idea of global optimal matching. But the Hungarian algorithm they adopted can only work in a one-to-one assignment manner. So far, for the CNN based detectors in one-to-many scenarios, a global optimal assigning strategy remains uncharted.

最近的许多工作试图使标签分配程序更具适应性，目的是进一步提高检测性能。GuidedAnchoring[40]没有使用预定义的锚，而是基于无锚机制生成锚，以更好地适应各种物体的分布。MetaAnchor[44]提出了一个锚点生成函数，从任意定制的先验盒中学习动态锚点。NoisyAnchors[19]提出了基于分类和定位损失的软标签和锚点重配机制。FreeAnchor[48]基于IoU为每个gt构建top-k锚点候选，然后提出一个检测定制的可能性，在每个候选集内进行正/负划分。ATSS[47]提出了一种自适应的样本选择策略，采用每个gt的最接近的锚集合的IoU值的平均值+std作为正/负阈值。PAA [16] 假设正负样本的联合损失分布遵循高斯分布。因此，它使用GMM来拟合正负样本的分布，然后使用正样本分布的中心作为正/负划分的边界。AutoAssign[51]通过自动确定空间和尺度维度上的正/负值，以完全数据驱动的方式处理标签分配。这些方法探讨了单个对象的最佳分配策略，而没有从全局角度考虑背景信息。DeTR[3]研究了全局最优匹配的思想。但是他们采用的匈牙利算法只能以一对一的分配方式工作。到目前为止，对于基于CNN的检测器在一对多的情况下，全局最优分配策略仍然是未知的。

3. Method

In this section, we first revisit the definition of the Optimal Transport problem and then demonstrate how we formulate the label assignment in object detection into an OT problem. We also introduce two advanced designs which we suggest adopting to make the best use of OTA.

在这一节中，我们首先重温了最优传输问题的定义，然后演示了我们如何将物体检测中的标签分配制定为一个OT问题。我们还介绍了两个先进的设计，我们建议采用这两个设计来最好地利用OTA。

3.1. Optimal Transport

The Optimal Transport (OT) describes the following problem: supposing there are m suppliers and n demanders in a certain area. The i -th supplier holds s_i units of goods while the j -th demander needs d_j units of goods. Transporting cost for each unit of good from supplier i to demander j is denoted by c_{ij} . The goal of OT problem is to find a transportation plan $\pi^* = \{\pi_{ij} \mid i$

$= 1, 2, \dots, m, j = 1, 2, \dots, n\}$, according to which all goods from suppliers can be transported to demanders at a minimal transportation cost:

最佳运输 (OT) 描述了以下问题: 假设在某一地区有 m 个供应商和 n 个需求者。第 i 个供应商持有 s_i 个单位的货物, 而第 j 个需求者需要 d_j 个单位的货物。每个单位的货物从供应商 i 到需求者 j 的运输成本用 c_{ij} 表示。OT 问题的目标是找到一个运输计划 $\pi^* = \{\pi_{ij} \mid i = 1, 2, \dots, m, j = 1, 2, \dots, n\}$, 根据这个计划, 供应商的所有货物都能以最小的运输成本运送给需求者。

$$\begin{aligned}
& \min_{\pi} \quad \sum_{i=1}^m \sum_{j=1}^n c_{ij} \pi_{ij} \\
& \text{s.t.} \quad \sum_{i=1}^m \pi_{ij} = d_j, \quad \sum_{j=1}^n \pi_{ij} = s_i, \\
& \quad \sum_{i=1}^m s_i = \sum_{j=1}^n d_j, \\
& \quad \pi_{ij} \geq 0, \quad i = 1, 2, \dots, m, j = 1, 2, \dots, n.
\end{aligned} \tag{1}$$

This is a linear program which can be solved in polynomial time. In our case, however, the resulting linear program is large, involving the square of feature dimensions with anchors in all scales. We thus address this issue by a fast iterative solution, named Sinkhorn-Knopp [5] (described in Appendix.)

这是一个线性程序, 可以在多项式时间内得到解决。然而, 在我们的案例中, 产生的线性程序很大, 涉及到所有尺度的锚的特征维度的平方。因此, 我们通过快速迭代解决这个问题, 命名为 Sinkhorn-Knopp [5] (在附录中描述。)

3.2. OT for Label Assignment

In the context of object detection, supposing there are m gt targets and n anchors (across all FPN [20] levels) for an input image I , we view each gt as a supplier who holds k units of positive labels (i.e., $s_i = k, i = 1, 2, \dots, m$), and each anchor as a demander who needs one unit of label (i.e., $d_j = 1, j = 1, 2, \dots, n$). The cost c_{ij}^{fg} for transporting one unit of positive label from gt_i to anchor a_j is defined as the weighted summation of their cls and reg losses:

在物体检测的背景下, 假设一个输入图像 I 有 m 个 gt 目标和 n 个锚点 (跨越所有的 FPN [20] 级别), 我们把每个 gt 看作是持有 k 个单位正标签的供应商 (即 $s_i = k, i = 1, 2, \dots, m$), 而每个锚点是需要一个单位标签的需求者 (即 $d_j = 1, j = 1, 2, \dots, n$)。从 gt_i 向锚框 a_j 运输一单位正标签的成本 c_{ij}^{fg} 被定义为它们的 cls 和 reg 损失的加权总和。

$$\begin{aligned}
c_{ij}^{fg} = & L_{cls}(P_j^{cls}(\theta), G_i^{cls}) + \\
& \alpha L_{reg}(P_j^{box}(\theta), G_i^{box}),
\end{aligned} \tag{2}$$

where θ stands for model's parameters. P_j^{cls} and P_j^{box} denote predicted cls score and bounding box for a_j . G_i^{cls} and G_i^{box} denote ground truth class and bounding box for gt_i . L_{cls} and L_{reg} stand for cross entropy loss and IoU Loss [46]. One can also replace these two losses with Focal Loss [21] and GIoU [34]/SmoothL1 Loss [11]. α is the balanced coefficient.

其中 θ 代表模型的参数。 P^{cls}_j 和 P^{box}_j 表示预测的cls分数和 a_j 的边界盒。 G^{cls}_i 和 G^{box}_i 表示地面真实类和 gt_i 的边界盒。 L_{cls} 和 L_{reg} 代表交叉熵损失和IoU损失[46]。也可以用Focal Loss[21]和GloU[34]/SmoothL1 Loss[11]来代替这两种损失。 α 是平衡系数。

Besides positive assigning, a large set of anchors are treated as negative samples during training. As the optimal transportation involves all anchors, we introduce another supplier – background, who only provides negative labels. In a standard OT problem, the total supply must be equal to the total demand. We thus set the number of negative labels that background can supply as

$n - m \times k$. The cost for transporting one unit of negative label from background to a_j is defined as:

除了正向匹配，在训练过程中，大量的锚框被当作负面样本。由于最优运输涉及到所有的锚，我们引入了另一个供应商--背景，他只提供负面标签。在一个标准的OT问题中，**总供给必须等于总需求**。因此，我们将背景可以提供的负面标签的数量设定为

$n-m \times k$ 。将一个单位的负面标签从背景运输到 a_j 的成本被定义为。

$$c_j^{bg} = L_{cls}(P_j^{cls}(\theta), \emptyset), \quad (3)$$

where \emptyset means the background class. Concatenating this $c^{bg} \in R^{(1 \times n)}$ to the last row of $c^{fg} \in R^{(m \times n)}$, we can get the complete form of the cost matrix $c \in R^{(m+1) \times n}$. The supplying vector s should be correspondingly updated as:

其中 \emptyset 指的是背景类。将这个 $c^{bg} \in R^{(1 \times n)}$ 与 $c^{fg} \in R^{(m \times n)}$ 的最后一行串联起来，我们可以得到成本矩阵 $c \in R^{(m+1) \times n}$ 的完整形式。供给向量 s 应相应地更新为。

$$s_i = \begin{cases} k, & \text{if } i \leq m \\ n - m \times k, & \text{if } i = m + 1. \end{cases} \quad (4)$$

As we already have the cost matrix c , supplying vector $s \in R^{(m+1)}$ and demanding vector $d \in R^n$, the optimal transportation plan $\pi^* \in R^{(m+1) \times n}$ can be obtained by solving this OT problem via the off-the-shelf Sinkhorn-Knopp Iteration [5]. After getting π^* , one can decode the corresponding label assigning solution by assigning each anchor to the supplier who transports the largest amount of labels to them. The subsequent processes (e.g., calculating losses based on assigning result, back-propagation) are exactly the same as in FCOS [38] and ATSS [47]. Noted that the optimization process of OT problem only contains some matrix multiplications which can be accelerated by GPU devices, hence OTA only increases the total training time by less than 20% and is totally cost-free in testing phase.

由于我们已经有了成本矩阵 c 、供应向量 $s \in R^{(m+1)}$ 和需求向量 $d \in R^n$ ，最优运输计划 $\pi^* \in R^{(m+1) \times n}$ 可以通过现成的Sinkhorn-Knopp迭代法[5]解决这个OT问题而得到。在得到 π^* 后，人们可以通过将每个锚分配给向其运输最大数量的标签的供应商来解码相应的标签分配方案。随后的过程（例如，根据分配结果计算损失，反向传播）与FCOS[38]和ATSS[47]中完全相同。注意到OT问题的优化过程只包含一些矩阵乘法，可以通过GPU设备加速，因此OTA只增加了不到20%的总训练时间，在测试阶段完全没有成本。

3.3. Advanced Designs Center Prior.

Previous works [47, 16, 48] only select positive anchors from the center region of objects with limited areas, called Center Prior. This is because they suffer from either a large number of ambiguous anchors or poor statistics in the subsequent process. Instead of relying on statistical characteristics, our OTA is based on global optimization methodology and thus is naturally resistant to these two issues. Theoretically, OTA can assign any anchor within the region of gts' boxes as a positive sample. However, for general detection datasets like COCO, we find the Center Prior still benefit the training of OTA. Forcing detectors focus on potential positive areas (i.e., center areas) can help stabilize the training process, especially in the early stage of training, which will lead to a better final performance. Hence, we impose a Center Prior to the cost matrix. For each gt, we select r 2 closest anchors from each FPN level according to the center distance between anchors and gts 2 . As for anchors not in the r 2 closest list, their corresponding entries in the cost matrix c will be subject to an additional constant cost to reduce the possibility they are assigned as positive samples during the training stage. In Sec. 4, we will demonstrate that although OTA adopts a certain degree of Center Prior like other works [38, 47, 48] do, OTA consistently outperforms counterparts by a large margin when r is set to a large value (i.e., large number of potential positive anchors as well as more ambiguous anchors).

以前的工作[47, 16, 48]只从面积有限的物体的中心区域选择积极的锚点，称为中心优先。这是因为他们在后续过程中要么存在大量的模糊锚点，要么存在糟糕的统计数据。我们的OTA不依赖于统计学特征，而是基于全局优化方法，因此自然能抵抗这两个问题。理论上，OTA可以将gts的盒子区域内的任何锚指定为阳性样本。然而，对于像COCO这样的一般检测数据集，我们发现Center Prior仍然有利于OTA的训练。迫使检测器专注于潜在的积极区域（即中心区域）可以帮助稳定训练过程，特别是在训练的早期阶段，这将导致更好的最终性能。因此，我们对成本矩阵施加一个中心先验。对于每个gt，我们根据锚和gt之间的中心距离，从每个FPN级别中选择 r 个最近的锚。对于不在 r 最接近列表中的锚，它们在成本矩阵 c 中的相应条目将受到一个额外的恒定成本的影响，以减少它们在训练阶段被分配为正样本的可能性。在第4节中，我们将证明，尽管OTA像其他作品[38, 47, 48]一样采用了一定程度的Center Prior，但当 r 被设定为大值时（即大量潜在的正向锚以及更多模糊的锚），OTA始终以较大的优势胜过同行。

Dynamic k Estimation.

Intuitively, the appropriate number of positive anchors for each gt (i.e., s_i in Sec. 3.1) should be different and based on many factors like objects' sizes, scales, and occlusion conditions, etc. As it is hard to directly model a mapping function from these factors to the positive anchor's number, we propose a simple but effective method to roughly estimate the appropriate number of positive anchors for each gt based on the IoU values between predicted bounding boxes and gts. Specifically, for each gt, we select the top q predictions according to IoU values. These IoU values are summed up to represent this gt's estimated number of positive anchors. We name this method as Dynamic k Estimation. Such an estimation method is based on the following intuition: The appropriate number of positive anchors for a certain gt should be positively correlated with the number of anchors that well-regress this gt. In Sec. 4, we present a detailed comparison between the fixed k and Dynamic k Estimation strategies. A toy visualization of OTA is shown in Fig. 2. We also describe the OTA's completed procedure including Center Prior and Dynamic k Estimation in Algorithm 1.

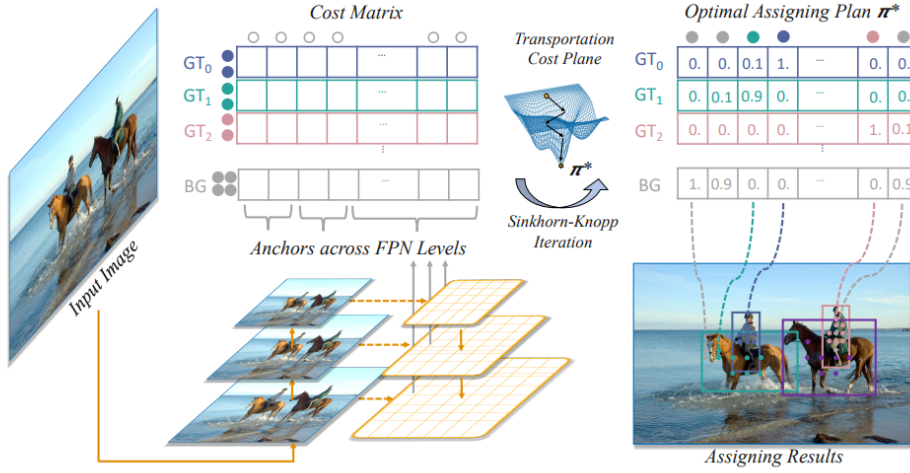


Figure 2. An illustration of Optimal Transport Assignment. *Cost Matrix* is composed of the pair-wise *cls* and *reg* losses between each anchor-gt pair. The goal of finding the best label assigning is converted to solve the best transporting plan which transports the labels from suppliers (i.e. GT and BG) to demanders (i.e. anchors) at a minimal transportation cost via Sinkhorn-Knopp Iteration.

直观地讲，每个gt（即第3.1节中的 s_i ）的合适的正锚数量应该是不同的，并基于许多因素，如物体的大小、比例和遮挡条件等。由于很难直接建立一个从这些因素到正锚数量的映射函数，我们提出了一个简单而有效的方法，根据预测的边界框和gt之间的IoU值来粗略估计每个gt的适当的正锚数量。具体来说，对于每个gt，我们根据IoU值选择前 q 个预测值。这些IoU值的总和代表了该gt的估计正锚的数量。我们将这种方法命名为动态 k 估计。这样的估计方法是基于以下的直觉。对某一指标而言，适当的正向锚点数量应该与能很好地还原这一指标的锚点数量呈正相关关系。在第4节中，我们将详细比较固定 k 和动态 k 估计策略。图2显示了OTA的一个玩具可视化。我们还描述了OTA的完成过程，包括算法1中的Center Prior和动态 k 估计。

Algorithm 1 Optimal Transport Assignment (OTA)

Input:

I is an input image

A is a set of anchors

G is the *gt* annotations for objects in image I

γ is the regularization intensity in Sinkhorn-Knopp Iter.

T is the number of iterations in Sinkhorn-Knopp Iter.

α is the balanced coefficient in Eq. 2

Output:

π^* is the optimal assigning plan

- 1: $m \leftarrow |G|, n \leftarrow |A|$
 - 2: $P^{\text{cls}}, P^{\text{box}} \leftarrow \text{Forward}(I, A)$
 - 3: $s_i (i = 1, 2, \dots, m) \leftarrow \text{Dynamic } k \text{ Estimation}$
 - 4: $s_{m+1} \leftarrow n - \sum_{i=1}^m s_i$
 - 5: $d_j (j = 1, 2, \dots, n) \leftarrow \text{OnesInit}$
 - 6: pairwise *cls* cost: $c_{\text{cls}}^{ij} = \text{FocalLoss}(P_j^{\text{cls}}, G_i^{\text{cls}})$
 - 7: pairwise *reg* cost: $c_{\text{reg}}^{ij} = \text{IoULoss}(P_j^{\text{box}}, G_i^{\text{box}})$
 - 8: pairwise Center Prior cost: $c_{ij}^{\text{cp}} \leftarrow (A_j, G_i^{\text{box}})$
 - 9: *bg cls* cost: $c_{\text{cls}}^{\text{bg}} = \text{FocalLoss}(P_j^{\text{cls}}, \emptyset)$
 - 10: *fg* cost: $c^{\text{fg}} = c_{\text{cls}} + \alpha c_{\text{reg}} + c_{\text{cp}}$
 - 11: compute final cost matrix c via concatenating $c_{\text{cls}}^{\text{bg}}$ to the last row of c^{fg}
 - 12: $v^0, u^0 \leftarrow \text{OnesInit}$
 - 13: **for** $t=0$ **to** T **do**:
 - 14: $u^{t+1}, v^{t+1} \leftarrow \text{SinkhornIter}(c, u^t, v^t, s, d)$
 - 15: compute optimal assigning plan π^* according to Eq. ??
 - 16: **return** π^*
-