

How to Backpropagate through Hungarian in Your DETR?

Lingji Chen, Alok Sharma, Chinmay Shirore, Chengjie Zhang, Balarama Raju Buddharaju
 Motional
 Boston, MA, USA
 motional.com

Abstract

The DEtection TRansformer (DETR) approach, which uses a transformer encoder-decoder architecture and a set-based global loss, has become a building block in many transformer based applications. However, as originally presented, the assignment cost and the global loss are not aligned, i.e., reducing the former is likely but not guaranteed to reduce the latter. And the issue of gradient is ignored when a combinatorial solver such as Hungarian is used. In this paper we show that the global loss can be expressed as the sum of an assignment-independent term, and an assignment-dependent term which can be used to define the assignment cost matrix. Recent results on generalized gradients of optimal assignment cost with respect to parameters of an assignment problem are then used to define generalized gradients of the loss with respect to network parameters, and backpropagation is carried out properly. Our experiments using the same loss weights show interesting convergence properties and a potential for further performance improvements.

1. Introduction

The seminal DEtection TRansformer (DETR) paper [5] views object detection as a direct set prediction problem. Its main ingredients are a set-based global loss that forces unique predictions via bipartite matching, and a transformer encoder-decoder architecture. DETR, together with its later variants such as Deformable DETR [31], has become a building block for many transformer based approaches, e.g. [3, 4, 18, 19, 21, 25, 28, 30] in the area of detection and tracking, to name just a few.

In [5], the network produces a set of predictions whose number is larger than the number of ground truth boxes. An assignment problem is defined and solved by the Hungarian algorithm, and the optimally matched boxes are used to define a global loss to be minimized through backpropagation. The loss accounts for three aspects of the prediction: (1) The classes of the matched boxes should be those of their assigned ground truth, (2) the positions and sizes of these boxes should be their assigned ground truth, and (3) the classes of the non-matched boxes should be background. Intuitively, the criterion according to which the matching is made, i.e., the total assignment cost, should be aligned with the global loss, such that reducing the former necessarily reduces the latter.

But this is not the case in [5], because the matching cost is defined differently from the loss. The terms accounting for classes use raw probability instead of cross entropy; the heuristic reason given is relative

scale to the loss from geometry. One may argue that the scaling problem is the same as in the loss, and can be explicitly dealt with using an additional scaling hyperparameter. The above terms are also limited to only the matched boxes; the reason given is not very convincing, because each prediction has a different probability of being background and therefore the missing sum is not matching-independent.

In a larger context, we may view the Hungarian solver as just another module performing some operations; in this case, some discrete optimization. One may argue that it does not matter how the cost matrix is defined, *as long as we can account for the gradient properly*. Unfortunately this is not done in most DETR approaches. In the code released with the paper [5], the issue of gradient is ignored by surrounding the `matcher` code with `torch.no_grad()`, *i.e.*, gradient tracing is turned off when Hungarian solver is involved.

In this paper we show that the matter is actually quite simple. First, the global loss can be expressed as a sum of two terms: One is defined by the probabilities of being background of *all predictions*, regardless of matching. The other can be treated as the optimal cost of matching, *if the cost matrix is suitably defined*.

Thus, to perform backpropagation on the global loss, we only need to answer the question: What is the gradient of the optimal cost with respect to the parameters defining the assignment problem? Fortunately this question, or rather a more general question involving Integer Linear Programming, has been answered in [11].

We have applied the new solution to the same training problem as in [5], and observed that the new solution produces desirable results, *i.e.*, loss decreases, convergence shows, and the model performance improves as training progresses. Note that very limited experiments have been conducted so far and the authors will be performing more extensive experimentation in future work.

This paper is organized as follows. In Sec. 2 we introduce the reader to the larger problem of differentiability when the network contains an optimization module. In Sec. 3 we present our new solution that not only aligns the matching cost with the global loss, but also simplifies the implementation. Experiments for the same training problem as in [5] with the proposed modifications are presented in Sec. 4 to show the efficacy of the new algorithm, and Sec. 5 summarizes the key takeaways. We plan to release our code.

2. Related Work

2.1. The general framework of Deep Declarative Networks

In [13] a new framework with end-to-end learnable models is presented, where a *declarative node* is defined by its desired behavior, in contrast to a conventional, *imperative node* that is defined by an explicit forward function. The desired behavior is typically stated as an objective and constraints in a mathematical optimization problem, and backpropagation can be performed, when the domain is continuous, with the help of the implicit function theorem, without having to know the actual solver used within. Previous work that inspired this generalization includes [2] that proposes an approach to differentiating through disciplined convex programs, and [12] that presents results on differentiating parameterized argmin and argmax problems in a lower-level which binds variables in the objective of an upper-level problem.

Since in a DETR based approach, a module solves an assignment problem, it can be considered a declarative node. However, because it involves discrete optimization, the gradient is not obtained through implicit function theorem, but through some other means, which will be described in the following

sections.

2.2. Differentiating combinatorial solvers by approximating the piecewise constant solution function

An interesting, general approach is proposed in [26] to approximate the piecewise constant solution function with an informative, differentiable function for a discrete optimization problem with a linear objective function. More specifically, to use the assignment problem as an example, let x be the input data together with the network weight (at the current iteration), I be the indices selected by the optimal solution, C be the optimal weight, and L be the loss to be minimized. If we start with $\partial L/\partial I$, then we would need $\partial I/\partial x$ to complete the backpropagation. However, the solution function from x to I is piecewise constant because there are only a finite number of values for I , and therefore its gradient is mostly zero or undefined. This is overcome in [26] by approximating the function with one that is differentiable and informative with non-zero gradients.

We would argue that it is often more convenient to specify the loss in terms of the optimal value (rather than the optimal solution), and start with $\partial L/\partial C$. Thus the utility of [26] is limited.

2.3. Sinkhorn operator based approaches for matching

To solve an assignment problem is to find a permutation matrix, where each row and each column has one and only one entry being “1” while the rest being “0.” This can be cast as a learning problem where the 1’s are the truth labels and the network output is interpreted as the probabilities of the entries. However, these probabilities have a structural constraint, *i.e.*, the matrix they form should be a doubly stochastic matrix (DSM) whose rows and columns all sum up to one. A bi-level optimization problem can be formulated where the lower level finds the best DSM approximation to the network output, which can be obtained through an iterative procedure called Sinkhorn-Knopp algorithm [24]. The gradient can be computed by unrolling the iterations which are differentiable. This is essentially the approach in [8, 20, 22] and many recent works.

Since unrolling the iterations takes a lot of memory, [10] gives a closed form expression of the gradient in the form of solutions to linear equations, based on the Karush–Kuhn–Tucker (KKT) conditions at the optimal value presumably achieved at the end of the iterations.

2.4. Differentiable assignment solver through relaxation and projection

In [29] the Integer Linear Programming (ILP) assignment problem is relaxed to a Linear Programming (LP) problem, and Dykstra’s algorithm [9] is used to project the solution into feasible space. This makes the solver differentiable.

2.5. Network approximated assignment solvers

To obtain differentiability, an approach is to approximate a combinatorial solver with a neural network whose operations are differentiable. This is achieved in [15, 17, 23, 27] and others. In [27], non-differentiable tracking metrics are also approximated by neural networks.

3. Our New Solution

3.1. The global loss, rearranged

We adopt some of the notations in [5] but also define new ones. Let $y = \{y_j\}_{j=1}^M$ be the set of M ground truth objects, and $\hat{y} = \{\hat{y}_i\}_{i=1}^N$ the set of N predictions. At inference time, thresholding probabilities or scores will give us a subset of the N predictions as actual objects. At training time,

we want to select exactly M out of N predictions to correspond to the ground truth. Since typically $N \gg M$, we want to solve an assignment problem with a *rectangular* cost matrix (as is done in the *code* released by [5], not the *paper* itself), not a *square* one. Thus, instead of talking about a permutation, we define an injective mapping $s : G \rightarrow B$ from the ground truth index set $G \triangleq \{1, 2, \dots, M\}$ to the predicted box index set $B \triangleq \{1, 2, \dots, N\}$. Define the matched set $B_1 \triangleq s(G)$, then the unmatched set is $B_2 \triangleq B \setminus B_1$. The inverse mapping $s^{-1} : B_1 \rightarrow G$ is naturally defined. In other words, starting from the ground truth indexed by j , we find the assigned prediction indexed by $s(j)$. And starting from a prediction indexed by i and already known to have had a ground truth assigned to it, we find the ground truth index $s^{-1}(i)$.

As in [5], let $y_j = (c_j, b_j)$ be the ground truth with object class label c_j and bounding box vector b_j . Let $\mathcal{L}_{\text{box}}(b_j, \hat{b}_{s(j)})$ be the loss of predicted bounding box $\hat{b}_{s(j)}$ when it is considered to represent b_j . Let $\hat{p}_i(c_{s^{-1}(i)})$ denote the probability of the target object class in the i -th prediction, and $\hat{p}_i(\emptyset)$, of background.

For a given set of network weight vector w , and a given mapping s , the loss has to do with three parts:

1. the object class probabilities in the matched set B_1 ,
2. the background probabilities in “the rest” set B_2 , and
3. the losses \mathcal{L}_{box} between G and B_1 .

The third part is the same as in [5] and details will be omitted in this paper. The first two parts depend on the likelihood

$$L \triangleq \prod_{j \in B_1} \hat{p}_j(c_{s^{-1}(j)}) \prod_{k \in B_2} \hat{p}_k(\emptyset). \quad (1)$$

Note that both terms depend on the mapping s . However, the following is a term that does not depend on the mapping s :

$$L_\emptyset \triangleq \prod_{i \in B} \hat{p}_i(\emptyset). \quad (2)$$

Thus the ratio

$$\frac{L}{L_\emptyset} = \prod_{j \in B_1} \frac{\hat{p}_j(c_{s^{-1}(j)})}{\hat{p}_j(\emptyset)} \quad (3)$$

only has to do with the index sets G and B_1 . The global loss (which is called Hungarian loss in [5]) is therefore

$$\begin{aligned} \mathcal{L}(y, \hat{y}, s) &= -\log\left(\frac{L}{L_\emptyset}\right) - \log(L_\emptyset) + \sum_{j=1}^M \mathcal{L}_{\text{box}}(b_j, \hat{b}_{s(j)}) \\ &= \sum_{j=1}^M \left(-\log \hat{p}_{s(j)}(c_j) + \log \hat{p}_{s(j)}(\emptyset) + \mathcal{L}_{\text{box}}(b_j, \hat{b}_{s(j)}) \right) - \sum_{i=1}^N \log \hat{p}_i(\emptyset). \end{aligned} \quad (4)$$

3.2. The new cost matrix

In the above loss, we will try to write the matching dependent terms as a total assignment cost, *i.e.*,

$$\begin{aligned} \mathcal{L}(y, \hat{y}, s) &= \sum_{j=1}^M \left(-\log \hat{p}_{s(j)}(c_j) + \log \hat{p}_{s(j)}(\emptyset) + \mathcal{L}_{\text{box}}(b_j, \hat{b}_{s(j)}) \right) - \sum_{i=1}^N \log \hat{p}_i(\emptyset) \\ &\triangleq \sum_{j=1}^M \mathcal{C}(s(j), j) + \mathcal{L}_{\text{background}}, \end{aligned} \tag{5}$$

where we define the (i, j) -th entry of an $N \times M$ cost matrix \mathcal{C} as

$$\mathcal{C}(i, j) \triangleq -\log \hat{p}_i(c_j) + \log \hat{p}_i(\emptyset) + \mathcal{L}_{\text{box}}(b_j, \hat{b}_i). \tag{6}$$

It is now clear that to minimize the global loss, we take two steps:

1. Solve an assignment problem using the cost matrix defined in Eq. (6):

$$\hat{s} = \operatorname{argmin}_{s: G \rightarrow B} \sum_{j=1}^M \mathcal{C}(s(j), j). \tag{7}$$

2. Minimize

$$\mathcal{L}(y, \hat{y}, \hat{s}) \triangleq \mathcal{L}_{\text{assign}} + \mathcal{L}_{\text{background}}. \tag{8}$$

3.3. The generalized gradients

The generalized gradients of the optimal assignment cost $\mathcal{L}_{\text{assign}}$ with respect to the network weights w can be obtained using the solution given in [11]. Thus the generalized gradients of $\mathcal{L}(y, \hat{y}, \hat{s})$ can be obtained and used in training [6].

More specifically, the cost matrix \mathcal{C} in Eq. (6) is defined by the ground truth y and the network predictions $\hat{y}(w)$ as a function of the network weight vector w . If we stack up the columns of \mathcal{C} to form a cost vector $c(w)$, then the assignment problem can be formulated as an Integer Linear Programming (ILP) problem with parameters $(c(w), A, b)$ where both A and b are constants specifying the inequality constraints of a valid assignment [7]. – We can pad the cost matrix with almost zero random values to form a square one with equality constraints, so that the results in [11] can be applied.

Let \hat{u} be an optimal solution (in the form of a column vector) to the assignment problem; practically speaking it is almost always unique. Then according to Algorithm 1 in [11], we have

$$\frac{\partial \mathcal{L}_{\text{assign}}}{\partial w} = \hat{u}^T \frac{\partial c}{\partial w}. \tag{9}$$

This gives us a way to calculate the gradients with `pytorch`: We treat the optimal assignment solution \hat{u} as a constant; this means that we take only the *numerical values* of the cost matrix and call a Hungarian solver, under `torch.no_grad()`. Then we can sum up the gradient-attached optimal cost according to the solution \hat{u} , together with the background terms, and then backpropagate the loss in Eq. (8).

3.4. The differences

We belabor the point by noting the following differences of our new solution compared to [5]:

- The assignment cost and the global loss are now aligned, and this is conceptually satisfying.
- The unmatched predictions have different probabilities of background in them that depend on the assignment, and therefore we do not ignore them in the cost matrix.
- The cross entropy loss is used in our cost matrix as opposed to the raw probabilities; the problem with scaling, if it still exists, can be solved separately.
- The challenge of the rectangular cost matrix is solved by using the ratio in Eq. (3) that results in the cost matrix in Eq. (6), *i.e.*, every row subtracts the cross entropy loss that corresponds to the background.
- There is a normalization by the number of boxes in their global loss definition for GIOU based losses, which is not done in their total assignment cost. The loss as expressed by Eq. (4) does not justify the normalization. After all, a batch having more boxes *should* have a bigger impact.

4. Experiments

We perform quantitative evaluation of the proposed global loss and cost matrix formulation using DETR architecture [5] on the COCO dataset [16]. For a comparison of performance, we first train DETR as is, and then train the model again using the proposed Assignment Aligned DETR as described in the previous section.

4.1. Dataset and Metrics

We use COCO2017 detection dataset [16], containing 118k training and 5k validation images. The images are annotated with bounding boxes and segmentation but for our purposes only the bounding boxes are relevant. For comparison with the DETR methodology, we report the integral AP metric used by COCO. This AP is averaged over multiple Intersection over Union (IoU) values. We also present plots showing how COCO AP metrics evolve over training epochs.

4.2. Training Details

For our experiments, we train DETR using the source code provided by its authors on [github](#). For the baseline using bi-partite matching loss as well as the proposed global loss, the model is trained as is using DETR setup hyper-parameters except that we use a 100 epoch training schedule, batch size of 2 per GPU and with auxiliary losses disabled. Note that DETR was trained using a 300 epoch training schedule with learning rate drops.

4.3. Experiment Results

Results of evaluation on COCO2017 validation set after 100 epochs are presented in Table 1. We also present plots showing how loss values changed over the 100 epoch training cycles in Fig 1a along with plots showing how the AP changed during the experiment.

The loss plots in Fig 1a indicate that with the new formulation, the loss consistently decreased and exhibited signs of convergence. The AP plots in Fig 1b, Fig 1c & Fig 1d also show that the network is

consistently able to learn the classes from data and its performance increased as the training progressed. If compared to the baseline DETR, the plots indicate that with the new loss formulation, the network performance did not exceed the baseline in the limited amount of experiments performed. However, it should be noted that the overall trend of metrics shows that conceptually speaking, the proposed mathematical formulation where the matching cost aligns with the global loss, produces desirable results specifically as loss decreases the network performance increases.

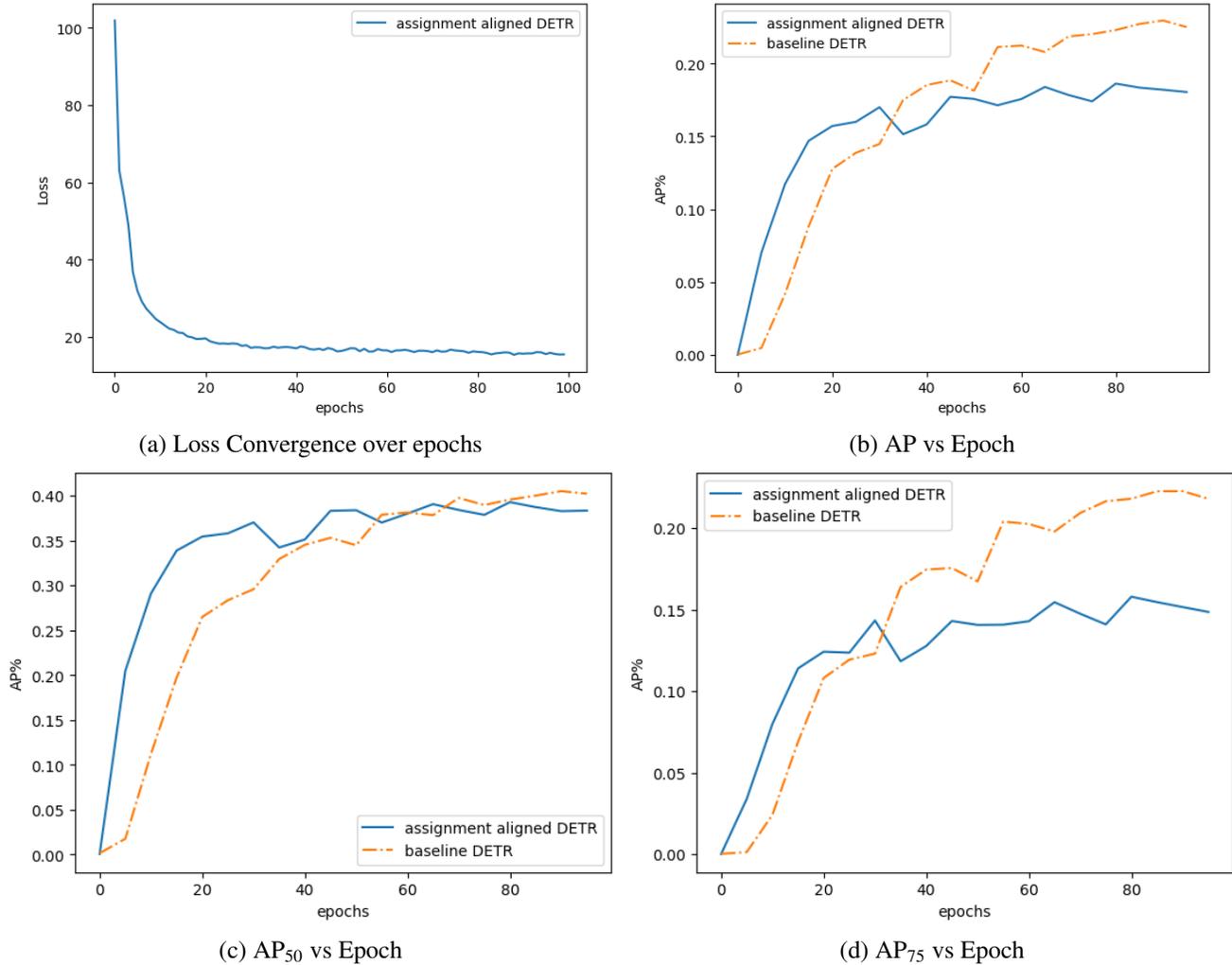


Figure 1. Results of training DETR with proposed loss for 100 epoch cycle

Model	Matching Strategy	Loss Minimized	AP	AP ₅₀	AP ₇₅	AP _S	AP _M	AP _L
DETR -R50	Hungarian	Set-Based Loss	24.1	42.40	235.5	8.5	25.1	39.0
DETR -R50	Hungarian	Assignment Cost	18.90	39.5	15.9	4.10	18.7	34.9

Table 1. Table shows the results of evaluation of DETR trained using loss defined in [5] and proposed loss implementation on COCO validation set .

5. Conclusion

We have presented a new formulation of assignment cost that can be applied to DETR like architectures which use Hungarian matching. This cost is aligned with the global loss calculation, and leads to proper use of the gradient even with the presence of a Hungarian solver. With the limited experiments conducted using DETR code-base, it can be concluded that the formulation is valid and produces desirable results. Given the promising results with initial experiments, the authors expect that matching SOTA results with a hyper-parameter search will be feasible, and such experiments will be part of the authors' future work. The authors also expect that the theoretically sound approach will encourage the research community to develop further and produce better performing models.

References

- [1] Sharath Adavanne, Archontis Politis, and Tuomas Virtanen. Differentiable tracking-based training of deep learning sound source localizers. pages 211–215, 2021.
- [2] Akshay Agrawal, Brandon Amos, Facebook Ai, Shane Barratt, Stephen Boyd, Steven Diamond, and J Zico Kolter. Differentiable convex optimization layers. *Advances in Neural Information Processing Systems*, 32, 2019. 2
- [3] Xuyang Bai, Zeyu Hu, Xinge Zhu, Qingqiu Huang, Yilun Chen, Hongbo Fu, and Chiew-Lan Tai. Transfusion: Robust lidar-camera fusion for 3d object detection with transformers. pages 1090–1099, 2022. 1
- [4] Jiarui Cai, Mingze Xu, Wei Li, Yuanjun Xiong, Wei Xia, Zhuowen Tu, and Stefano Soatto. Memot: Multi-object tracking with memory. pages 8090–8100, 2022. 1
- [5] Nicolas Carion, Francisco Massa, Gabriel Synnaeve, Nicolas Usunier, Alexander Kirillov, and Sergey Zagoruyko. End-to-end object detection with transformers. volume 12346 LNCS, pages 213–229. Springer Science and Business Media Deutschland GmbH, 2020. 1, 2, 3, 4, 6, 7
- [6] Frank H Clarke. *Optimization and nonsmooth analysis*. SIAM, 1990. 5
- [7] David F Crouse. On implementing 2d rectangular assignment algorithms. *IEEE Transactions on Aerospace and Electronic Systems*, 52:1679–1696, 2016. 5
- [8] Rodrigo Santa Cruz, Basura Fernando, Anoop Cherian, and Stephen Gould. Visual permutation learning. *IEEE transactions on pattern analysis and machine intelligence*, 41:3100–3114, 2018. 3
- [9] Richard L Dykstra. An algorithm for restricted least squares regression. *Journal of the American Statistical Association*, 78:837–842, 1983. 3
- [10] Marvin Eisenberger, Aysim Toker, Laura Leal-Taixé, Florian Bernard, and Daniel Cremers. A unified framework for implicit sinkhorn differentiation. pages 509–518, 2022. 3
- [11] Xi Gao, Han Zhang, Aliakbar Panahi, and Tom Arodz. Combinatorial losses through generalized gradients of integer linear programs. *arXiv preprint arXiv:1910.08211*, 2019. 2, 5
- [12] Stephen Gould, Basura Fernando, Anoop Cherian, Peter Anderson, Rodrigo Santa Cruz, and Edison Guo. On differentiating parameterized argmin and argmax problems with application to bi-level optimization. *arXiv preprint arXiv:1607.05447*, 2016. 2
- [13] Stephen Gould, Richard Hartley, and Dylan Campbell. Deep declarative networks: A new hope. *arXiv preprint arXiv:1909.04866*, 2019. 2
- [14] Georg Hess and William Ljungbergh. Transforming the field of multi-object tracking, 2021.
- [15] Mengyuan Lee, Yuanhao Xiong, Guanding Yu, and Geoffrey Ye Li. Deep neural networks for linear sum assignment problems. *IEEE Wireless Communications Letters*, 7:962–965, 2018. 3

- [16] Tsung-Yi Lin, Michael Maire, Serge J. Belongie, Lubomir D. Bourdev, Ross B. Girshick, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C. Lawrence Zitnick. Microsoft COCO: common objects in context. *CoRR*, abs/1405.0312, 2014. 6
- [17] He Liu, Tao Wang, Congyan Lang, Songhe Feng, Yi Jin, and Yidong Li. Glan: A graph-based linear assignment network. *arXiv preprint arXiv:2201.02057*, 2022. 3
- [18] Fan Ma, Mike Zheng Shou, Linchao Zhu, Haoqi Fan, Yilei Xu, Yi Yang, and Zhicheng Yan. Unified transformer tracker for object tracking. pages 8781–8790, 2022. 1
- [19] Tim Meinhardt, Alexander Kirillov, Laura Leal-Taixe, and Christoph Feichtenhofer. Trackformer: Multi-object tracking with transformers. pages 8844–8854, 2022. 1
- [20] Gonzalo Mena, David Belanger, Scott Linderman, and Jasper Snoek. Learning latent permutations with gumbel-sinkhorn networks. *arXiv preprint arXiv:1802.08665*, 2018. 3
- [21] Ishan Misra, Rohit Girdhar, and Armand Joulin. An end-to-end transformer model for 3d object detection. pages 2906–2917, 2021. 1
- [22] Ioannis Papakis, Abhijit Sarkar, and Anuj Karpatne. Gcnmatch: Graph convolutional neural networks for multi-object tracking via sinkhorn normalization. *arXiv preprint arXiv:2010.00067*, 2020. 3
- [23] Athena Psalta, Vasileios Tsironis, and Konstantinos Karantzas. Transformer-based assignment decision network for multiple object tracking. *arXiv preprint arXiv:2208.03571*, 2022. 3
- [24] Richard Sinkhorn and Paul Knopp. Concerning nonnegative matrices and doubly stochastic matrices. *Pacific Journal of Mathematics*, 21:343–348, 1967. 3
- [25] Peize Sun, Jinkun Cao, Yi Jiang, Rufeng Zhang, Enze Xie, Zehuan Yuan, Changhu Wang, and Ping Luo. Transtrack: Multiple object tracking with transformer. *arXiv preprint arXiv:2012.15460*, 2020. 1
- [26] Marin Vlastelica, Anselm Paulus, Vit Musil, Georg Martius, and Michal Rolinek. Differentiation of blackbox combinatorial solvers. 2019. 3
- [27] Yihong Xu, Aljosa Osep, Yutong Ban, Radu Horaud, Laura Leal-Taixé, and Xavier Alameda-Pineda. How to train your deep multi-object tracker. pages 6787–6796, 2020. 3
- [28] Fangao Zeng, Bin Dong, Tiancai Wang, Xiangyu Zhang, and Yichen Wei. Motr: End-to-end multiple-object tracking with transformer. *arXiv preprint arXiv:2105.03247*, 2021. 1
- [29] Xiaohui Zeng, Renjie Liao, Li Gu, Yuwen Xiong, Sanja Fidler, and Raquel Urtasun. Dmm-net: Differentiable mask-matching network for video object segmentation. pages 3929–3938, 2019. 3
- [30] Xingyi Zhou, Tianwei Yin, Vladlen Koltun, and Philipp Krähenbühl. Global tracking transformers. pages 8771–8780, 2022. 1
- [31] Xizhou Zhu, Weijie Su, Lewei Lu, Bin Li, Xiaogang Wang, and Jifeng Dai. Deformable detr: Deformable transformers for end-to-end object detection. *arXiv preprint arXiv:2010.04159*, 2020. 1