

Fast and Flexible Human Pose Estimation with HyperPose

Yixiao Guo*
Peking University
China

Jiawei Liu*
Tongji University
China

Guo Li*
Imperial College London
United Kingdom

Luo Mai
University of Edinburgh
United Kingdom

Hao Dong
Peking University
China

ABSTRACT

Estimating human pose is an important yet challenging task in multimedia applications. Existing pose estimation libraries target reproducing standard pose estimation algorithms. When it comes to customising these algorithms for real-world applications, none of the existing libraries can offer both the flexibility of developing custom pose estimation algorithms and the high-performance of executing these algorithms on commodity devices. In this paper, we introduce HyperPose, a novel flexible and high-performance pose estimation library. HyperPose provides expressive Python APIs that enable developers to easily customise pose estimation algorithms for their applications. It further provides a model inference engine highly optimised for real-time pose estimation. This engine can dynamically dispatch carefully designed pose estimation tasks to CPUs and GPUs, thus automatically achieving high utilisation of hardware resources irrespective of deployment environments. Extensive evaluation results show that HyperPose can achieve up to 3.1x~7.3x higher pose estimation throughput compared to state-of-the-art pose estimation libraries without compromising estimation accuracy. By 2021, HyperPose has received over 1000 stars on GitHub and attracted users from both industry and academy.

CCS CONCEPTS

• **Computing methodologies** → **Computer vision**; • **Software and its engineering** → **Software performance**.

KEYWORDS

Pose Estimation, Computer Vision, High-performance Computing

ACM Reference Format:

Yixiao Guo, Jiawei Liu, Guo Li, Luo Mai, and Hao Dong. 2021. Fast and Flexible Human Pose Estimation with HyperPose. In *Proceedings of the 29th ACM International Conference on Multimedia (MM '21)*, October 20–24, 2021, Virtual Event, China. ACM, New York, NY, USA, 4 pages. <https://doi.org/10.1145/3474085.3478325>

1 INTRODUCTION

Multimedia applications, such as interactive gaming, augmented reality and self-driving cars, can greatly benefit from accurate and fast human pose estimation. State-Of-The-Art (SOTA) pose estimation

*Equal contributions. Yixiao implemented the model development library. Jiawei implemented the model inference engine. Guo prototyped the initial version of HyperPose.

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.

MM '21, October 20–24, 2021, Virtual Event, China

© 2021 Association for Computing Machinery.

ACM ISBN 978-1-4503-8651-7/21/10...\$15.00

<https://doi.org/10.1145/3474085.3478325>

algorithms (e.g., OpenPose [1], PoseProposal [14], and PifPaf [10]) pre-process video streams, use neural networks to infer human anatomical key points, and then estimate the human pose topology.

In practice, achieving both high accuracy and high-performance in pose estimation is however difficult. On the one hand, achieving high-accuracy pose estimation requires users to deeply customise standard pose estimation algorithms (e.g., OpenPose and PifPaf), so that these algorithms can accurately reflect the characteristics of user-specific deployment environments (e.g., object size, illumination, number of humans), thus achieving high accuracy. On the other hand, deeply customised pose estimation algorithms can contain various *non-standard* computational operators in pre-processing and post-processing data, making such algorithms difficult to always exhibit high performance in using commodity embedded platforms (e.g., NVIDIA Jetson and Google TPU Edge)

When using existing pose estimation libraries to develop custom applications, users often report several challenges. High-performance C++-based libraries such as OpenPose [1] and AlphaPose [3] focus on specific pose estimation algorithms. They do not provide intuitive APIs for users to customise pose estimation algorithms based on the requirements of their specific deployment environments. These libraries are also optimised for certain hardware platforms. When re-targeting them to new hardware platforms, users must largely modify their internal execution runtime, which is non-trivial for most pose estimation algorithm developers. Furthermore, users could also use high-level pose estimation libraries such as TF-Pose [8] and PyTorch-Pose [18]. These libraries offer users with Python APIs to easily declare various pose estimation algorithms. However, this easiness comes with a large performance overhead, making these libraries incapable of handling real-world deployment where high-resolution images are ingested at high speed [17].

In this paper, we introduce HyperPose, a flexible and fast library for human pose estimation. The design and implementation of HyperPose makes the following contributions:

(i) Flexible APIs for developing custom pose estimation algorithms. HyperPose provides flexible Python APIs for developing custom pose estimation algorithms. These APIs consist of those for customising the pipelines of pose estimation algorithms, the architectures of deep neural networks, training datasets, training hyper-parameters [11], data pre-processing pipelines, data post-processing pipelines, and the strategy of paralleling neural network training on GPUs. We show that, using these APIs, users can declare a wide range of commonly used pose estimation algorithms while customising these algorithms for high estimation accuracy.

(ii) High-performance model inference engine for executing custom pose estimation algorithms. HyperPose can achieve high-performance in executing custom pose estimation algorithms. This is achieved through a novel *high-performance algorithm execution engine*. This engine is designed as a streaming dataflow [13].

This dataflow can take custom computational operators for implementing custom pose estimation logic. These operators can be dynamically dispatched onto parallel CPUs and GPUs, thus keeping computational resources always busy, irrespective of model architectures and hardware platforms. The implementations of these operators are also highly optimised, mainly by carefully leveraging the optimised computer-vision library: OpenCV, and the high-performance model inference library: TensorRT.

We study both the API easiness and the performance of HyperPose. The API study shows that HyperPose can provide better flexibility via building and customising pose estimation algorithms. The test-bed experiments further show that HyperPose can outperform the state-of-the-art optimised pose estimation framework: OpenPose by up to 3.1x in terms of the processing throughput of high-resolution images.

2 DESIGN AND IMPLEMENTATION

In this section, we present the design principles and implementation details of HyperPose.

2.1 Expressive Programming APIs

HyperPose aims to support different types of users in developing pose estimation algorithms. There are users who would like to find suitable algorithms and adapt them for their applications. To support this, HyperPose allows users to customise the pipeline of a typical pose estimation algorithm. Other users would like to further modify the components in a pose estimation pipeline. For example, they often need to control how a deep neural network is being trained, and the data pre-processing/post-processing operators being used. To support them, HyperPose allow users to plug in user-defined components in a pose estimation pipeline.

2.1.1 Algorithm development APIs. HyperPose provides a set of high-level APIs¹ to relieve users from the burden of assembling the complex pose estimation system. The APIs are in three modules, including *Config*, *Model*, and *Dataset*. *Config* exposures APIs to configure the pose estimation system, while *Model* and *Dataset* offer APIs to construct the concrete model, dataset, and the development pipeline. With each API, users can configure the architecture for different algorithms (e.g., OpenPose, PoseProposal), backbone networks (e.g., MobileNet [7] and ResNet[6]), and the training or evaluation datasets (e.g., COCO and MPII). For the development procedure, users can configure the hyper-parameters (e.g., learning rate and batch size), the distributed training via the KungFu library [12] option (e.g., using a single or multiple GPUs), the training strategy (e.g., adopting pre-training and adaptation stage), and the format of storing a trained model for further deployment (e.g., ONNX and TensorRT UFF). The rich configuration options enable users to efficiently adapt the off-the-shelf models.

2.1.2 Algorithm customisation APIs. HyperPose users can flexibly customise pose estimation algorithms to best fit in with their specific usage scenarios. This is achieved by providing common interfaces for key components in the algorithms. For example, to implement custom neural networks, users could inherit from the

Model class defined in HyperPose, and warping the custom computation logic into the corresponding member functions. As long as use the self-defined model to replace the preset model options during configuration, the custom model is enabled. The same practice applies to enabling a custom dataset. These customised components are then automatically integrated into the pose estimation system by HyperPose. The *Model* module further exposures processing modules including preprocessor, postprocessor and visualizer, which allow users to assemble their own development pipeline. By doing this, HyperPose makes its APIs flexible to support extensive customisation of its pose estimation algorithms.

2.2 High-performance Execution Engine

Designing a high-performance execution engine for human pose estimation is challenging. A human pose estimation pipeline consists of video stream ingesting, pre-processing (e.g., resizing and data layout switching), GPU model inference, CPU post-processing, and result exportation (e.g., visualisation).

These computation steps must be collaboratively completed using heterogeneous devices: CPUs, GPUs, and I/O devices (e.g., disk, cameras, etc.). To achieve real-time inference, the engine must maximise the efficiency of using *all* these devices through parallelism. We make two key designs in our inference engine:

2.2.1 Streaming pose estimation dataflow. We abstract the computation steps shared by most pose estimation algorithms and implement these steps as *dataflow operators*. The operators can be asynchronously executed in a streaming dataflow. The topology of the dataflow is implemented to be static for provisioning better optimisation, thus maximising the processing throughput.

Figure 1 illustrates the dataflow implemented in HyperPose. The source for the dataflow is often a video stream produced by a real-time device (e.g., camera). The dataflow ingests the images into a *decoding operator* (see ❶). The decoded images are resized and transformed to an expected data layout (e.g., channel-first) so that they can be fed into a neural network (see ❷). This neural network is executed by an *inference operator* (see ❸) on GPU. This operator loads the checkpoint of a neural network trained by the HyperPose Python platform or other supported libraries. The computed activation map from the neural network is given to a *post-processing operator* for parsing human pose topology (see ❹). The topology is placed onto the original image at a *visualisation operator* (see ❺).

Each operator pair shares a concurrent FIFO channel, manipulated by CPU threads. The operators only block when there is no incoming record and sleep until they are notified, thanks to the condition variable mechanism. Benefiting from this, HyperPose can ensure the pose estimation dataflow fully utilise parallel heterogeneous processors including CPUs and GPUs. Besides, HyperPose fully masks I/O latency (e.g., waiting for images from a camera) by overlapping I/O operations with computation operations.

2.2.2 Hybrid dataflow operator scheduler. We design a *hybrid dataflow operator scheduler* in the execution engine (see ❻ in Figure 1), for further improving the CPU/GPU utilisation in addition to the pipeline parallelism. Regarding GPU utilisation, we leverage dynamic batching in front of the inference operator, encouraging the inference operator to take a larger batch as input each time. In concert with

¹<https://hyperpose.readthedocs.io/en/latest/>

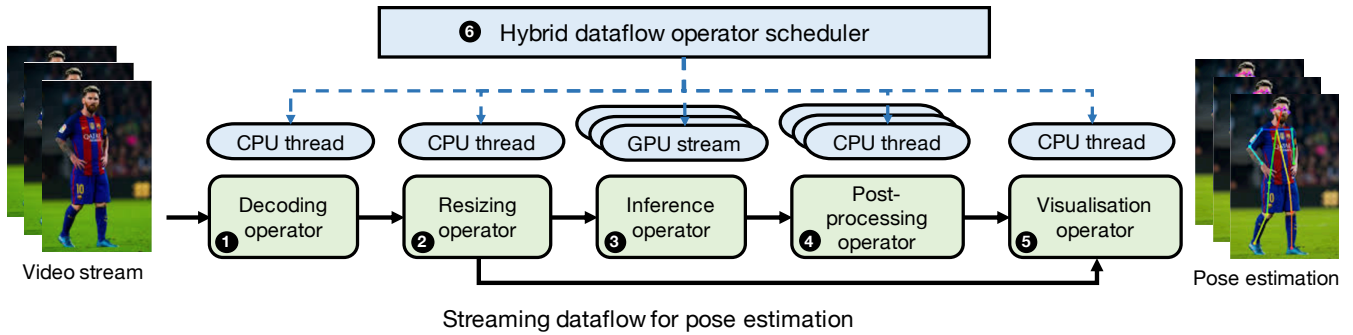


Figure 1: Architecture of the HyperPose C++ execution engine.

the streaming dataflow mechanism, the batching slot only accumulates more input tensors when the GPU becomes the bottleneck. Such optimisation is beneficial since i) batching reduces the times of GPU kernel launch thus improving GPU processing throughput, ii) when the GPU is the bottleneck, batching gets enhanced to alleviate the congestion [9], and iii) when the GPU is not the bottleneck, batching gets weakened for not deteriorating the per-image response delay. As to CPU threads scheduling, we implemented an asynchronous thread-level communication mechanism based on conditional variables. When the bottleneck happens in a CPU-based operator (e.g., the post-processing operator), the working threads of non-blocking operators will fall asleep to save CPU cycles for the bottleneck until the next round starts.

2.3 Implementation and Compatibility

HyperPose supports 3 classes of pose estimation algorithms: 1) PAF[1] (i.e., OpenPose), 2) PoseProposal[14], and 3) PifPaf[10].

Since developing and deploying pose estimation algorithms have different objectives, HyperPose separates the implementation for training and inference but makes their ecosystem well-compatible.

The training library is a Python library implemented using Tensorflow and TensorLayer [2] for DNN construction with Numpy and other common libraries for post-processing.

For maximum performance, the whole inference engine is notably implemented in C++ 17 for massive low-level code optimisation and parallelism. The GPU inference operator is based on NVidia TensorRT, one of the fastest DNN inference libraries. The imaging-related operations are based on OpenCV, and the dataflow scheduler is implemented by the C++ standard thread library.

The implementation of HyperPose is compatible with many pose estimation algorithms, such as DEKR [4] and CenterNet [19]. These algorithms share the bottom-up architectures as those (e.g., OpenPose and PifPaf) implemented in HyperPose. They are thus easy to be implemented as extensions to HyperPose.

3 EVALUATION

Our evaluation of HyperPose is driven by two questions: (i) How flexible is its API when developing and customising real-world pose estimation algorithms? (ii) How fast is its execution engine in practical deployment environments?

²For fair accuracy comparison, the weights are from PifPaf and PoseProposal libraries.

Table 1: API study. * denotes single-human datasets only.

	Algorithm	Dataset	DNN	Config.	Ext.
TF-Pose	1	1	5	5	✗
PyTorch-Pose	2	3*	13	26	✗
HyperPose	3	2	10	30	✓

Table 2: Performance Evaluation of Inference Engine. ²

Configuration	Baseline FPS	Our FPS (Operators)	Our FPS (Scheduler)
OpenPose (VGG19)	8[5]	19.78	27.32
OpenPose (MobileNet)	8.5[8]	50.89	84.32
LWOpenPose (ResNet50)	N/A	38.09	63.52
LWOpenPose (TinyVGG)	N/A	66.62	124.92
PoseProposal (ResNet18)	47.6[16]	212.42	349.17
PifPaf (ResNet50)	14.8[15]	18.5	44.13

Table 3: Accuracy Evaluation of Development Platform.

Configuration	Original Accuracy(map)	Our Accuracy(map)
OpenPose (VGG19)	58.4	57.0
OpenPose (MobileNet)	28.1	44.2
LWOpenPose (MobileNet)	42.8	46.1
LWOpenPose (Resnet50)	N/A	48.2
LWOpenPose (TinyVGG)	N/A	47.3

3.1 API Study

We compare our API design with existing Python pose estimation libraries: TF-Pose and PyTorch-Pose. Other libraries such as OpenPose and AlphaPose are dedicated algorithm implementations without customised extensions, we thus exclude them here.

In Table 1, our comparison follows five metrics: (i) the number of pre-defined pose estimation algorithms, (ii) the number of pre-defined datasets, (iii) the number of pre-defined backbone deep

neural networks (DNNs), (iv) the total number of pre-defined configurations of the pose estimation system, and (v) the ability to extend the library to support custom algorithms.

TF-Pose only supports 1 algorithm, 1 dataset, and 5 DNNs, resulting in up to 5 configurations in its design space. PyTorch-Pose is more flexible, which supports 2 pose estimation algorithms, 3 datasets, and 13 DNNs, summing up to 26 system configurations. However, PyTorch-Pose only covers single-human scenarios. This is attributed to the insufficient performance of its algorithm execution engine in multi-human scenarios. Contrastingly, HyperPose provides 3 algorithms, 2 datasets, and 10 DNNs, thus supporting up to 30 system configurations. Moreover, in all these Python libraries, HyperPose is the only one that supports the extension of new pose estimation algorithms and provides abstract processing modules that allow users to build their own development pipeline.

3.2 Performance Evaluation

Table 2 compares the existing libraries and HyperPose. All benchmarks are evaluated under the same configuration. The test-bed is of 6 CPU cores and 1 NVIDIA 1070Ti GPU. We measure the throughput of the pose estimation systems. The benchmark video stream comes from the *Crazy Uptown Funk Flashmob in Sydney* which contains 7458 frames with 640x360 resolution.

We first compare the performance of HyperPose with the OpenPose framework [1], which leverages Caffe as its backend and uses C++ for implementing pre-processing and post-processing. As is shown in Table 2, OpenPose is only able to achieve 8 FPS on 1070 Ti, which HyperPose can reach 27.32 FPS, outperforming the baseline by 3.1x. On one hand, this improvement is attributed to the careful use of the TensorRT library as the implementation of the inference operator. On the other hand, the hybrid dataflow operator scheduler makes the execution of HyperPose even 1.38x faster than the non-scheduled one. TF-Pose [8] leverages TensorFlow as its inference engine and its post-processing is implemented in C++ as well. When executing MobileNet-based OpenPose, it only achieves 8.5 FPS, which is 10x slower than HyperPose.

In addition to OpenPose-based algorithms, HyperPose also outperforms Chainer’s [16] implementation of Pose Proposal Network by 8 times. We verified the performance consistency by replacing the backbones and post-processing methods. For example, HyperPose also beats OpenPose when evaluating a smaller model (*i.e.*, MobileNet). This proves that the execution engine design is generic so that its benefits should be shared by all custom algorithms.

Table 3 shows the accuracy evaluation result of HyperPose in MS COCO dataset, which includes 5K images in val2017 and 40K images in val2014. The evaluation dataset split follows the original implementations.

4 CONCLUSION

When operating pose estimation in the wild, developers often find it challenging to customise pose estimation algorithms for high accuracy, while achieving real-time pose estimation using commodity CPUs and GPUs. This paper introduces HyperPose, a library for fast and flexible human pose estimation. HyperPose provides users with expressive APIs for declaring custom pose estimation algorithms. It also provides a high-performance model inference

engine that can efficiently utilise all parallel CPUs and GPUs. This engine enables HyperPose to achieve 3.1x better performance than existing libraries, while achieving the same accuracy in challenging pose estimation tasks.

ACKNOWLEDGMENTS

This project was supported by the National Natural Science Foundation of China-Youth Science Fund (No.62006006) and National Key R&D Program of China: New Generation Artificial Intelligence Open Source Community and Evaluation (No.2020AAA0103500), Topic: New Generation Artificial Intelligence Open Source Community Software and Hardware Infrastructure Support Platform (No.2020AAA0103501).

REFERENCES

- [1] Zhe Cao, Tomas Simon, Shih-En Wei, and Yaser Sheikh. 2017. Realtime multi-person 2d pose estimation using part affinity fields. In *Proceedings of the IEEE conference on computer vision and pattern recognition*. 7291–7299.
- [2] Hao Dong, Akara Supratak, Luo Mai, Fangle Liu, Axel Oehmichen, Simiao Yu, and Yike Guo. 2017. TensorLayer: A Versatile Library for Efficient Deep Learning Development. In *Proceedings of the 25th ACM International Conference on Multimedia (Mountain View, California, USA) (MM ’17)*. Association for Computing Machinery, New York, NY, USA, 1201–1204. <https://doi.org/10.1145/3123266.3129391>
- [3] Hao-Shu Fang, Shuqin Xie, Yu-Wing Tai, and Cewu Lu. 2017. Rmpe: Regional multi-person pose estimation. In *Proceedings of the IEEE international conference on computer vision*. 2334–2343.
- [4] Zigang Geng, Ke Sun, Bin Xiao, Zhaoxiang Zhang, and Jingdong Wang. 2021. Bottom-Up Human Pose Estimation via Disentangled Keypoint Regression. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*. 14676–14686.
- [5] Ginés Hidalgo 2020. *OpenPose: Real-time multi-person keypoint detection library*. <https://github.com/CMU-Perceptual-Computing-Lab/openpose>
- [6] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. 2016. Deep Residual Learning for Image Recognition. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*.
- [7] Andrew G Howard, Menglong Zhu, Bo Chen, Dmitry Kalenichenko, Weijun Wang, Tobias Weyand, Marco Andreetto, and Hartwig Adam. 2017. Mobilenets: Efficient convolutional neural networks for mobile vision applications. *arXiv preprint arXiv:1704.04861* (2017).
- [8] Dongwoo Kim Ildoo Kim. 2019. *tf-pose-estimation*. Retrieved June 7, 2021 from <https://github.com/ildoonet/tf-pose-estimation>
- [9] Alexandros Koliouisis, Pijika Watcharapichat, Matthias Weidlich, Luo Mai, Paolo Costa, and Peter Pietzuch. [n.d.]. CROSSBOW: Scaling Deep Learning with Small Batch Sizes on Multi-GPU Servers. *Proceedings of the VLDB Endowment* 12, 11 ([n. d.]).
- [10] Sven Kreiss, Lorenzo Bertoni, and Alexandre Alahi. 2019. Pifpaf: Composite fields for human pose estimation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 11977–11986.
- [11] Luo Mai, Alexandros Koliouisis, Guo Li, Andrei-Octavian Brabete, and Peter Pietzuch. 2019. Taming hyper-parameters in deep learning systems. *ACM SIGOPS Operating Systems Review* 53, 1 (2019), 52–58.
- [12] Luo Mai, Guo Li, Marcel Wagenländer, Konstantinos Fertakis, Andrei-Octavian Brabete, and Peter Pietzuch. 2020. KungFu: Making Training in Distributed Machine Learning Adaptive. In *14th USENIX Symposium on Operating Systems Design and Implementation (OSDI 20)*. 937–954.
- [13] Luo Mai, Kai Zeng, Rahul Potharaju, Le Xu, Steve Suh, Shivaram Venkataraman, Paolo Costa, Terry Kim, Saravanan Muthukrishnan, Vamsi Kuppala, et al. 2018. Chi: A scalable and programmable control plane for distributed stream processing systems. *Proceedings of the VLDB Endowment* 11, 10 (2018), 1303–1316.
- [14] Taiki Sekii. 2018. Pose proposal networks. In *Proceedings of the European Conference on Computer Vision (ECCV)*. 342–357.
- [15] Sven Kreiss 2021. *OpenPifPaf: Composite Fields for Semantic Keypoint Detection and Spatio-Temporal Association*. <https://github.com/openpifpaf/openpifpaf>
- [16] Satoshi Terasaki. 2018. *Chainer implementation of Pose Proposal Networks*. <https://github.com/Idain/chainer-pose-proposal-net>
- [17] Le Xu, Shivaram Venkataraman, Indranil Gupta, Luo Mai, and Rahul Potharaju. 2021. Move Fast and Meet Deadlines: Fine-grained Real-time Stream Processing with Cameo. In *18th {USENIX} Symposium on Networked Systems Design and Implementation ({NSDI} 21)*. 389–405.
- [18] Wei Yang. 2019. *pytorch-pose*. <https://github.com/bearpaw/pytorch-pose>
- [19] Xingyi Zhou, Dequan Wang, and Philipp Krähenbühl. 2019. Objects as points. *arXiv preprint arXiv:1904.07850* (2019).