# Deep Learning for Inertial Positioning: A Survey

Changhao Chen and Xianfei Pan

arXiv:2303.03757v2 [cs.RO] 20 Mar 2023

*Abstract*—Inertial sensors are widely utilized in smartphones, drones, robots, and IoT devices, playing a crucial role in enabling ubiquitous and reliable localization. Inertial sensor-based positioning is essential in various applications, including personal navigation, location-based security, and human-device interaction. However, low-cost MEMS inertial sensors' measurements are inevitably corrupted by various error sources, leading to unbounded drifts when integrated doubly in traditional inertial navigation algorithms, subjecting inertial positioning to the problem of error drifts. In recent years, with the rapid increase in sensor data and computational power, deep learning techniques have been developed, sparking significant research into addressing the problem of inertial positioning. Relevant literature in this field spans across mobile computing, robotics, and machine learning. In this article, we provide a comprehensive review of deep learning-based inertial positioning and its applications in tracking pedestrians, drones, vehicles, and robots. We connect efforts from different fields and discuss how deep learning can be applied to address issues such as sensor calibration, positioning error drift reduction, and multi-sensor fusion. This article aims to attract readers from various backgrounds, including researchers and practitioners interested in the potential of deep learning-based techniques to solve inertial positioning problems. Our review demonstrates the exciting possibilities that deep learning brings to the table and provides a roadmap for future research in this field.

*Index Terms*—Inertial Navigation, Deep Learning, Inertial Sensor Calibration, Pedestrian Dead Reckoning, Sensor Fusion, Visual-inertial Odometry

## I. INTRODUCTION

THE inertial Measurement Unit (IMU) is widely used in smartphones, drones, VR/AR devices, robotics, and Internet of Things (IoT) devices. It continuously measures linear velocity and angular rate and tracks the motion of these platforms, as illustrated in Figure 1. With the advancements in Micro-Electro-Mechanical Systems (MEMS) technology, today's MEMS IMUs are small, energy-efficient, and cost-effective. Inertial positioning (navigation) calculates attitude, velocity, and position based on inertial measurements, making it a crucial element in various location-based applications, including locating and navigating individuals in public places (e.g., universities, malls, airports), supporting security and safety services (e.g., aiding first-responders), enabling smart city/infrastructure, and facilitating human-device interaction. Compared to other positioning solutions such as vision or radio, inertial positioning is completely ego-centric, works indoors and outdoors, and is less affected by environmental
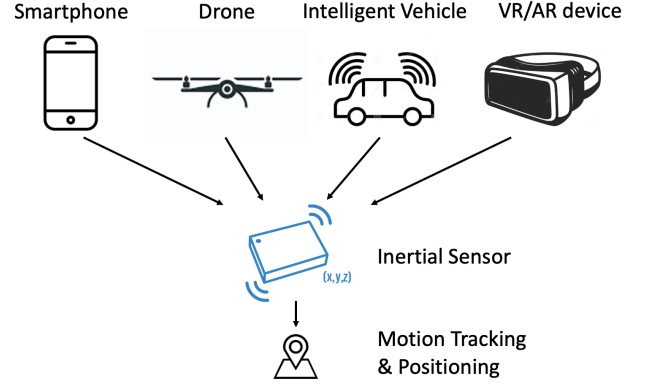
Fig. 1: Inertial sensors are ubiquitous in modern platforms such as smartphones, drones, intelligent vehicles, and VR/AR devices. They play a critical role in enabling completely ego-centric motion tracking and positioning, making them essential for a range of applications.

factors such as complex lighting conditions and scene dynamics.

Unfortunately, the measurements obtained from low-cost MEMS IMUs are subject to several error sources such as bias error, temperature-dependent error, random sensor noise, and random-walk noise. In classical inertial navigation mechanisms, angular rates are integrated into orientation, and based on the acquired attitude, acceleration measurements are transformed into the navigation frame. Finally, the transformed accelerations are doubly integrated into locations [1], [2]. Traditional inertial navigation algorithms are designed and described using concrete physical and mathematical rules. Under ideal conditions, sensor errors are small enough to allow hand-designed inertial navigation algorithms to produce accurate and reliable pose estimates. However, in real-world applications, inevitable measurement errors cause significant problems for inertial positioning systems without constraints, which can fail within seconds. In this process, even a minor error can be amplified exponentially, resulting in unbounded error drifts.

Previous researchers have attempted to address the problem of error drifts in inertial navigation by incorporating domain-specific knowledge or other sensor. In the context of pedestrian tracking, exploiting the periodicity of human walking is important, and the process of pedestrian dead reckoning (PDR) involves detecting steps, estimating step length and heading, and updating the user's location to mitigate error drifts from exponential to linear increase [3]. Zero-velocity update (ZUPT) involves attaching the IMU to the user's foot and detecting the zero-velocity phase, which is then used in Kalman filtering to correct inertial navigation states [4].

Platforms such as drones or robots equipped with other sensors such as cameras or LiDAR can significantly improve the performance of pure inertial solutions by effectively integrating inertial sensors with these modalities through filtering or smoothing [5]–[7]. However, these solutions have limitations in specific application domains and are unable to address the fundamental problem of inertial navigation.

Recently, deep learning has shown impressive performance in various fields, including computer vision, robotics, and signal processing [8]. It has also been introduced to address the challenges of inertial positioning. Deep neural network models have been leveraged to calibrate inertial sensor noises, reduce the drifts of inertial navigation mechanisms, and fuse inertial data with other sensor information. These research works have attracted significant attention, as they show potential for exploiting massive data to generate data-driven models instead of relying on concrete physical or mathematical models. With the rapid development of deep learning techniques, learning-based inertial solutions have become even more promising.

In this survey, we provide a comprehensive review of deep-learning-based approaches to inertial positioning, including measurement calibration, inertial positioning algorithms, and sensor fusion. We discuss the benefits and limitations of existing works and identify challenges and future opportunities in this research direction. Compared with other deep learning surveys, such as those focused on object detection [9], semantic segmentation [10], and robotics [11], *survey on deep learning based inertial positioning is relatively scarce and hard to find*. While a broader survey on machine learning enhanced inertial sensing does exist [12], our survey narrows the focus to deep learning based inertial positioning, providing deeper insights and analysis of the fast-evolving developments in this area over the past five years (2018-2022). Other relevant surveys, such as those focused on inertial pedestrian positioning [3], indoor positioning [13], step length estimation [14], and pedestrian dead reckoning [15], do not cover recent deep learning based solutions. *To the best of our knowledge, this article is the first survey that discusses deep learning based inertial positioning thoroughly and deeply.*

The rest of this survey is organized as follows. Section II briefly introduces classical inertial navigation mechanisms. Section III, IV and V survey deep learning based sensor calibration, inertial navigation algorithms, and sensor fusion. Section VII finally discusses the benefits, challenges and opportunities.

## II. Classical Inertial Navigation Mechanisms

This section provides an overview of classical inertial navigation mechanisms and highlights their limitations. It begins by presenting the inertial measurement model and classical strapdown inertial navigation method. Subsequently, two solutions that aim to reduce the drifts of inertial navigation system, namely pedestrian dead reckoning (PDR) and zero-velocity update (ZUPT), are discussed, with a specific focus on their applicability in pedestrian tracking scenarios. The section finally introduces sensor fusion approaches that integrate inertial data with information from other sensors.

### A. Inertial Measurement Model

Inertial measurements acquired from low-cost MEMS IMUs are often corrupted by various types of error sources, resulting in unbounded error drifts when integrated in strapdown inertial navigation systems (SINS). These error sources can be classified into two categories: deterministic errors and random errors [16]. Deterministic errors comprise bias error, non-orthogonality error, misalignment error, scale-factor error, and temperature-dependent error. On the other hand, random errors include random sensor noise and random-walk noise resulting from long-term operation, which are challenging to model and eliminate.

Raw IMU measurements, i.e. accelerations $\hat{\mathbf{a}}$ and angular rates $\hat{\boldsymbol{\omega}}$, can be formulated by

$$\hat{\mathbf{a}} = \mathbf{a} + \mathbf{b}_a + \mathbf{n}_a \tag{1}$$

$$\hat{\boldsymbol{\omega}} = \boldsymbol{\omega} + \mathbf{b}_\omega + \mathbf{n}_\omega \tag{2}$$

where $\mathbf{b}_a$ and $\mathbf{b}_\omega$ are acceleration bias and gyroscope bias, $\mathbf{n}_a$ and $\mathbf{n}_\omega$ are additive noises above accelerometer and gyroscope.

Traditionally, it is important to calibrate inertial sensors before running an inertial navigation algorithm that involves integrating inertial data into system states. One effective tool for achieving this is the Allan variance method [17], which models the random process of inertial sensor errors.

### B. Strapdown Inertial Navigation System

Inertial sensor measures linear accelerations $\mathbf{a}_b(t)$ and angular rates $\boldsymbol{\omega}_b^n(t)$ of attached user body at the timestep $t$. $b$ represents the body frame, while $n$ denotes the navigation (world) frame, i.e. the navigation frame. $\boldsymbol{\omega}_b^n(t)$ means that the angular rates of body frame with respect to the navigation frame. To simplify inertial motion model, this article assumes that the biases and noises of sensor in Equation 1 and 2 have been removed in the stage of inertial sensor calibration. $(\mathbf{R}, \mathbf{p})$ are defined orientation and position variables. From the kinematic model of IMU, we can have

$$\begin{cases} \mathbf{R}_b^n(t+1) = \mathbf{R}_b^n(t)\mathbf{R}_{b_{t+1}}^{b_t} \\ \mathbf{v}_n(t+1) = \mathbf{v}_n(t) + \mathbf{a}_n(t)dt \\ \mathbf{p}_n(t+1) = \mathbf{p}_n(t) + \mathbf{v}_n(t)dt + \frac{1}{2}\mathbf{a}_n(t)dt^2 \end{cases} \tag{3}$$

where $\mathbf{a}_n, \mathbf{v}_n, \mathbf{p}_n$ are acceleration, velocity and position in the navigation frame, $\mathbf{R}_b^n$ represents the rotation from the body frame to the navigation frame.

Firstly, orientation is updated by inferring the rotation matrix $\boldsymbol{\Omega}(t)$ via Rodriguez formula:

$$\begin{aligned} \boldsymbol{\Omega}(t) &= \mathbf{R}_{b_{t+1}}^{b_t} \\ &= \mathbf{I} + \sin(\boldsymbol{\sigma})\frac{[\boldsymbol{\sigma}\times]}{\boldsymbol{\sigma}} + (1 - \cos(\boldsymbol{\sigma}))\frac{[\boldsymbol{\sigma}\times]^2}{\boldsymbol{\sigma}^2}, \end{aligned} \tag{4}$$

where rotation vector $\boldsymbol{\sigma} = \boldsymbol{\omega}(t)dt$.

To update velocity, the accelerations in navigation frame can be expressed as a function of measured accelerations, i.e.

$$\mathbf{a}_n(t) = \mathbf{R}_b^n(t-1)\mathbf{a}_b(t) - \mathbf{g}_n \tag{5}$$

Then, the accelerations in navigation frame $\mathbf{a}_n(t)$ are integrated into the velocity in the navigation frame $\mathbf{v}_n(t)$, and the

location $\mathbf{p}_n(t)$ is finally updated by integrating the velocity via Equation 4.

As we can see, in this process, even a small measurement error can be exponentially amplified, leading to the problem of inertial error drifts. In the past, high-precision inertial sensors such as laser or fiber inertial sensors could keep the measurement error small enough to maintain the accuracy of INS. However, due to the size and cost limitations of current MEMS IMUs, compensation methods are necessary to mitigate the corresponding error drifts. One approach is to introduce domain-specific knowledge or other sensor information.

### C. Domain Specific Knowledge

*1) Pedestrian Dead Reckoning:* Pedestrian dead reckoning (PDR) is a method that leverages domain-specific knowledge about human walking to track pedestrian motion. PDR comprises three main steps: step detection, heading and stride length estimation, and location update [3]. In step detection, PDR uses the threshold of inertial data to identify step peaks or stances and segment the corresponding inertial data. Dynamic stride length estimation is then achieved via an empirical formula, known as the Weinberg formula [18], which considers the segmented accelerations and user's height. Similar to SINS, heading estimation is done by integrating gyroscope signals into orientation changes and adding orientation changes to the initial orientation to obtain the current heading. Finally, the estimated heading and stride length are used to update the pedestrian's location. By avoiding double integration of accelerations and incorporating a reliable stride estimation model, PDR effectively reduces inertial positioning drifts. However, inaccurate step detection and stride estimation can still occur, leading to large system error drifts. Moreover, PDR is limited to pedestrian navigation as it depends on the periodicity of human walking.

*2) Zero-velocity Update:* The Zero-velocity update (ZUPT) algorithm is designed to compensate for the errors of SINS by identifying the still phase of human walking and using zero-velocity as observations in a Kalman filter [4]. To facilitate the detection of the still phase, the IMU is typically attached to the user's foot, as it undergoes significant motion and reflects walking patterns well. Techniques such as peak-detection [19], zero-crossings [20], or auto-correlation [21] can be used to analyze the inertial data and segment the zero-velocity phase. Once the still phase is detected, zero-velocity is used as pseudo-measurements in the filtering process, thereby limiting the error drifts of open-loop integration. However, the effectiveness of ZUPT depends on the assumption that the user's foot remains completely still, and any incorrect still phase detection or small motion disturbances can cause navigation system drifts. Additionally, ZUPT is limited to pedestrian tracking.

### D. Integrating IMU with Other Sensors

Integrating the IMU with other sensors, such as camera [6], LiDAR [7], UWB [22], and magnetometer [23], can provide promising results as it allows for exploiting their
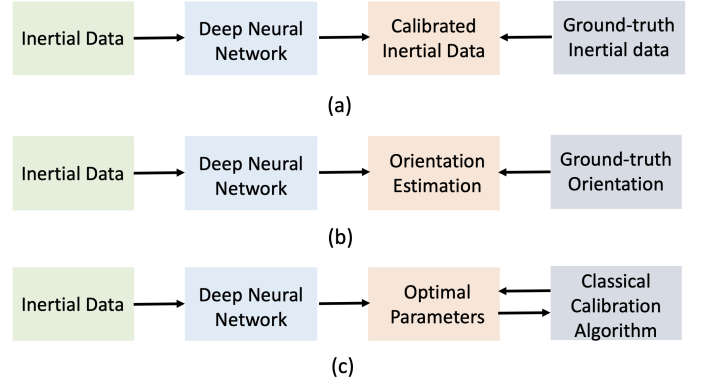


Fig. 2: An overview of existing deep learning based inertial sensor calibration methods

complementary properties. By fusing the data from multiple sensors, the accuracy and robustness of pose estimation can be significantly improved, making it a general solution for all platforms. However, in some scenarios, certain sensors, such as visual perception, may not be available or highly dependent on the environment, which can negatively affect the egocentric property of inertial positioning. Additionally, in sensor fusion approaches, it is essential to consider various factors such as sensor calibration, initialization, and time-synchronization.

### E. Discussion

As previously mentioned, classical inertial navigation methods are designed to solve specific problems within their respective domains. However, their performance is often limited due to real-world issues such as imperfect modeling, measurement errors, and environmental influences, resulting in inevitable error drifts. Researchers in the field of inertial navigation are therefore constantly searching for ways to build models that can tolerate measurement errors and mitigate system drifts. In addition to relying on Newtonian physical rules, it has been observed that domain-specific knowledge, whether it be an experienced human walking model or scene geometry, can serve as a useful constraint in reducing the error drifts of inertial positioning systems. One potential approach to improving inertial positioning accuracy and robustness is to exploit massive inertial data to extract domain-specific knowledge and construct a data-driven model. In the next sections, we will delve deeper into this problem and explore potential solutions.

### III. DEEP LEARNING BASED INERTIAL SENSOR CALIBRATION

Inertial measurements obtained from low-cost IMUs are often affected by various sources of noise, making it challenging to distinguish the true values from the sources of error. The error sources are a complex interplay of deterministic and random factors, further complicating the issue. To address the impact of measurement errors, the powerful nonlinear approximator capabilities of deep neural networks can be exploited. A natural approach is to develop a deep neural network that receives the raw inertial measurements as input

TABLE I: A summary of existing methods on deep learning based inertial sensor calibration.

| name | year | sensor | model | learning | target |
|------|------|--------|-------|----------|--------|
| Xiyuan et al. [24] | 2003 | gyro | 1-layer NN | SL | gyro drifts compensation |
| Chen et al. [25] | 2018 | gyro, acc | ConvNet | SL | inertial noise compensation |
| Esfahani et al. [26] | 2019 | gyro | LSTM | SL | gyroscope calibration |
| Nobre et al. [27] | 2019 | gyro, acc | Deep Q-Network | RL | optimal calibration parameters |
| Brossard et al. [28] | 2020 | gyro | ConvNet | SL | gyro corrections |
| Zhao et al. [29] | 2020 | gyro | LSTM | SL | gyroscope calibration |
| Huang et al. [30] | 2022 | gyro | Temporal ConvNet | SL | gyroscope calibration |
| Calib-Net [31] | 2022 | gyro | Dilated ConvNet | SL | gyroscope denoising |

- *Years* indicates the publication year of each work.
- *Sensors* indicates the sensors involved in each work. gyro and acc represent gyroscope and accelerometer respectively.
- *Model* indicates which module the framework consists of.
- *Learning* indicates how to train neural networks. SL and RL represent Supervised Learning and Reinforcement Learning.
- *Target* indicates what the model aims to solve or produce.

and produces the calibrated inertial measurements as output, representing the actual platform motion. By training this neural model on labeled datasets using stochastic gradient descent (SGD) [32], the inertial measurement errors can be implicitly learned and corrected by the neural network. It is important to note that the quality of the collected training dataset has a significant impact on the performance of the model.

Before the age of deep learning, attempts were made to use neural networks to learn the measurement errors of inertial sensors. For example, a 1-layer artificial neural network (ANN) [33] is proposed to model the distribution of gyro drifts, and is able to successfully approximate gyro drifts with such a 'shallow' network [24] . This method has an advantage over Kalman filtering (KF) based calibration methods in that it does not require setting hyper-parameters before use, such as the sensor noise matrix in KF.

In recent years, there has been increasing interest in using deep neural networks (DNN) with multiple layers to solve the inertial sensor calibration problem. With the addition of more layers, neural networks become more expressive and can learn complex relationships between the raw inertial measurements and the true motion of the vehicle. One approach, proposed by [25], uses a Convolutional Neural Network (ConvNet) to remove error noises from inertial measurements. They collected inertial data from two grades of IMU under given constant accelerations and angular rates. The ConvNet framework takes raw inertial measurements (from low-precision IMU) as inputs and tries to output acceleration and angular rate references (from high-precision IMU). Their experiment shows that deep learning can remove some of the sensor error and improve test accuracy. However, this work has not been validated in a real navigation setup, and thus it cannot demonstrate how learning-based sensor calibration reduces error drifts in inertial navigation. Both of the mentioned methods require reference data from high-precision IMUs as labels to train the networks, as shown in Figure 2 (a). However, acquiring reference data from high-precision IMUs can be costly.

In addition to directly learning from pseudo ground-truth IMU labels, another approach is to enable neural network-based calibration models to produce inertial data that can be integrated into more accurate orientation estimation. This is illustrated in Figure 2 (b). By producing more accurate
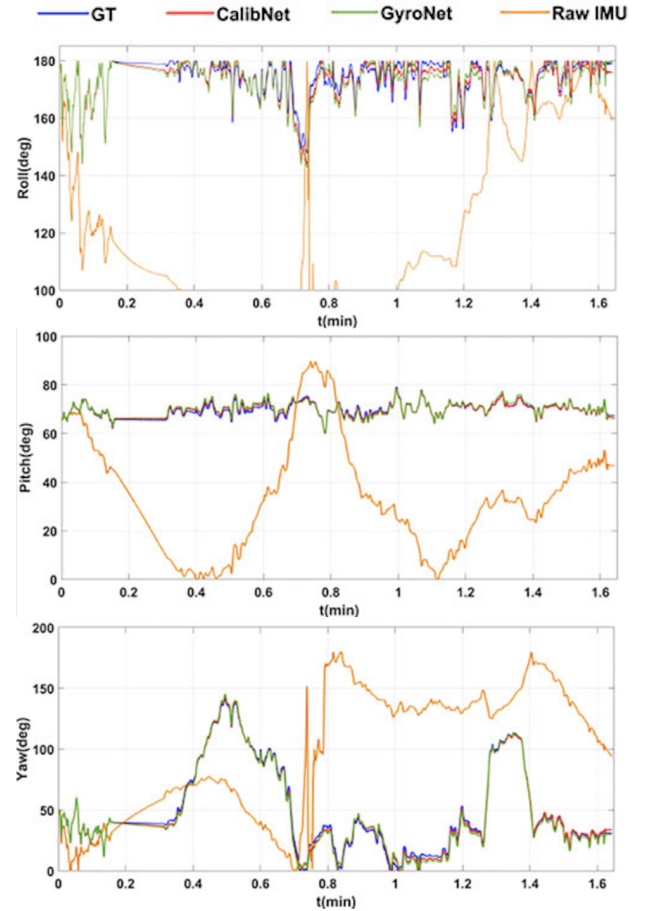


Fig. 3: An example of gyro calibration results (reprint from Calib-Net [31]). Compared with raw IMU integration, deep learning based calibration models significantly reduce attitude drifts.

orientation values, the neural network implicitly removes the corrupted noises above inertial data. For example, OriNet [26] inputs 3-dimensional gyroscope signals into an LSTM network [66] to obtain calibrated gyroscope signals, which are then integrated with the orientation at the previous timestep to generate orientation estimates at the current timestep. A loss function between orientation estimates and real orientation

TABLE II: A summary of existing methods on deep learning based inertial positioning.

| name | Year | Carrier | model | learning | target |
|---|---|---|---|---|---|
| IONet [34] | 2018 | Pedestrian, Trolley | LSTM | SL | location displacement |
| RIDI [35] | 2018 | Pedestrian | SVM, SVR | SL | velocity for inertial data calibration |
| Cortes et al. [36] | 2018 | Pedestrian | ConvNet | SL | velocity to constrain system drifts |
| Wagstaff et al. [37] | 2018 | Pedestrian | LSTM | SL | zero-velocity detection for ZUPT |
| Chen et al. [38] | 2019 | Pedestrian, Trolley | LSTM | TL | location displacement |
| AbolDeepIO [39] | 2019 | UAV | LSTM | SL | location displacement |
| RINS-W [40] | 2019 | Vehicle | RNN | SL | zero-velocity dection for KF |
| Feigl et al. [41] | 2019 | Pedestrian | LSTM | SL | walking velocity |
| Wang et al. [42] | 2019 | Pedestrian | LSTM | SL | walking heading for ZUPT |
| Yu et al. [43] | 2019 | Pedestrian | ConvNet | SL | adaptive zero-velocity detection |
| TLIO [44] | 2020 | Pedestrian | ConvNet | SL | 3D displacement and uncertainty for EKF |
| LIONet [45] | 2020 | Pedestrian | Dilated ConvNet | SL | lightweight inertial model |
| RoNIN [46] | 2020 | Pedestrian | LSTM, TCN | SL | velocity for inertial data calibration |
| Brossard et al. [47] | 2020 | Vehicle | ConvNet | SL | co-variance noise for KF |
| StepNet [48] | 2020 | Pedestrian | ConvNet, LSTM | SL | dynamic step length for PDR |
| Wang et al. [49] | 2020 | Pedestrian | ConvNet | SL | measurement noise for Kalman Filter |
| ARPDR [50] | 2020 | Pedestrian | TCN | SL | stride length and walking heading for PDR |
| IDOL [51] | 2021 | Pedestrian | LSTM | SL | device orientation and location |
| PDRNet [52] | 2021 | Pedestrian | ConvNet | SL | step length and heading for PDR |
| Buchanan et al. [53] | 2021 | Legged Robot | ConvNet | SL | integrate location displacement with leg odometry |
| Zhang et al. [54] | 2021 | Vehicle, UAV | RNN | SL | independent motion terms |
| Gong et al. [55] | 2021 | Pedestrian | LSTM | SL | fusing inertial data from two devices |
| NILoc [56] | 2022 | Pedestrian | ConvNet | SL | inertial relocalization |
| RIO [57] | 2022 | Pedestrian | DNN | UL | rotation-equivariance as supervision signal |
| Wang et al. [58] | 2022 | Pedestrian | DNN | SL | efficient and low-latent model |
| TinyOdom [59] | 2022 | Pedestrian, Vehicle | TCN+NAS | SL | deployment on resource-constrained device |
| CTIN [60] | 2022 | Pedestrian | Transformer | SL | velocity and trajectory prediction |
| DeepVIP [61] | 2022 | Vehicle | ConvNet, LSTM | SL | velocity and heading for car localization |
| Bo et al. [62] | 2022 | Pedestrian | ConvNet | TL | model-independent stride learning |
| OdoNet [63] | 2022 | Vehicle | ConvNet | SL | speed learning for ZUPT |
| A2DIO [64] | 2022 | Pedestrian | ConvNet, LSTM | SL | pose invariant odometry |
| LLIO [65] | 2022 | Pedestrian | MLP | SL | 3D displacement for lightweight odometry |

- *Year* indicates the publication year of each work.
- *Carrier* indicates the platform running inertial navigation.
- *Model* indicates which module the framework consists of.
- *Learning* indicates how to train neural networks. SL, TL and UL represent Supervised Learning, Transfer Learning and Unsupervised Learning.
- *Target* indicates what the model aims to solve or produce.

is defined and minimized for model training. OriNet has been evaluated on a public drone dataset, demonstrating an improvement in orientation performance of approximately 80%. A similar approach is [28], who calibrates gyroscope using ConvNet, reporting good attitude estimation accuracy. Calib-Net [31] is another ConvNet framework that denoises gyroscope data by extracting effective spatio-temporal features from inertial data. Calib-Net is based on dilation ConvNet [67] to compensate the gyro noise, as illustrated in Figure 3. This model is able to significantly reduce orientation error compared to raw IMU integration. When this learned inertial calibration model is incorporated into a visual-inertial odometry (VIO), it further improves localization performance and outperforms representative VIOs such as VINS-mono [6]. Other efforts in this direction include works by [29], [30].

Instead of directly calibrating inertial sensors with DNNs, some researchers have explored using DNNs to generate parameters that improve classical calibration algorithms, as shown in Figure 2 (c). One example is the work by [27], who models inertial sensor calibration as a Markov Decision Process and proposes to use deep reinforcement learning [68] to learn the optimal calibration parameters. The authors

demonstrated the effectiveness of their approach in calibrating inertial sensors for a visual-inertial odometry (VIO) system.

As discussed above, deep learning-aided inertial sensor calibration methods (listed in Table I) have shown promising results in removing corrupted sensor noises and improving the accuracy of inertial positioning systems. These methods do not require human intervention and can automatically learn error models. However, it is important to note that the learned error model is typically dependent on the specific sensor or platform used. Therefore, a change in sensor or user can result in different data distributions, leading to reduced performance of the learned model. Additionally, further analysis is needed to determine which types of noise can be effectively removed by learning-based calibration methods.

## IV. LEARNING TO CORRECT IMU INTEGRATION

In addition to sensor calibration, researchers are exploring various methods for using deep learning to construct inertial positioning models that can either partially or completely replace classical inertial navigation mechanisms. This section provides an overview of how deep learning can be used to correct IMU integration in general. Next sections will discuss
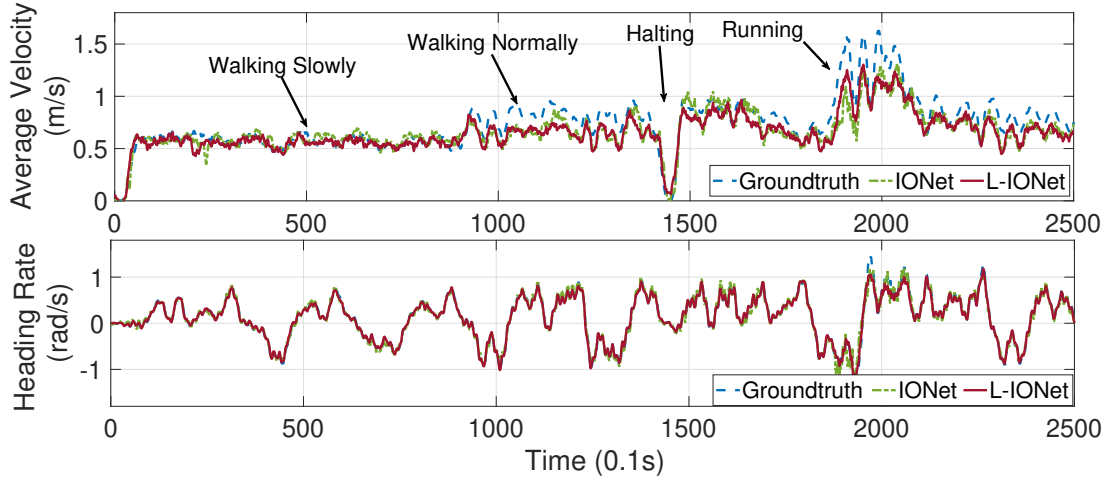
Fig. 4: The velocity of attached platform can be inferred from a sequence of inertial measurements via deep neural networks. (reprint from L-IONet [45])
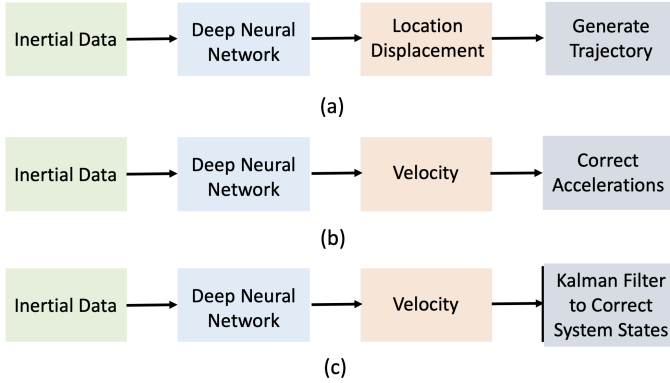


Fig. 5: An overview of existing methods on learning to correct IMU integration

deep learning approaches for pedestrian tracking applications, and present deep inertial solutions for vehicles, UAVs, and robots. A summary of existing works and their contributions is provided in Table II.

In deep learning-based inertial positioning approaches, a user's absolute velocity can be inferred from a sequence of IMU data using a deep neural network. This velocity information can then be used as a key constraint to reduce the drifts in IMU double integration. Figure 4 provides an example of velocity learning from IMU sequence, where the periodicity of human walking makes it easy to infer the user's moving velocity. Similar observations have been made for vehicles, UAVs, and robotic platforms, which will be discussed in Section VI. Existing works on applying learned velocity to correct IMU integration can generally be divided into three categories, as shown in Figure 5, and will be discussed as follows.

One category of deep learning models aims to learn location displacement, which is the average velocity multiplied by a fixed period of time, as illustrated in Figure 5(a). The approach proposed by [34] formulates inertial positioning as a sequential learning problem, where 2D motion displace-

ments in the polar coordinate, also known as polar vectors, are learned from independent windows of segmented inertial data. This is because the frequency of platform vibrations is relevant to the absolute moving speed, which can be measured by IMU, when tracking human or wheeled configurations. Based on this observation, they propose IONet, an LSTM-based framework for end-to-end learning of relative poses. Trajectories are generated by adding motion displacements together with initial locations. To train neural models, a large collection of data was collected from a smartphone-based IMU in a room with a high-precision visual motion tracking system (i.e., Vicon) to provide ground-truth pose labels. Once the model is trained, the IONet model can be used in areas outside the data-collection room. In a two-minute random pedestrian walking scenario, the localization error of IONet is within 3 meters 90% of the time, when evaluating across users, devices, and attachments, outperforming some classical PDR algorithms. In tracking trolley, IONet shows comparable performance over representative visual-inertial odometry and is even more robust in featureless areas. However, supervised learning-based IONet requires high-precision pose as training labels. When testing with data different from those in the training set, there will be performance degradation. To improve the generalization ability, [38] proposes MotionTransformer, which allows the inertial positioning model to self-adapt into new domains via generative adversarial network (GAN) [69] and domain adaptation [70], without the need for labels in new domains. To encourage more reliable inertial positioning, [71] is able to produce pose uncertainties along with poses, offering the belief in the extent to which the learned pose can be trusted. To allow full 3D localization, TLIO [44] proposes to learn 3D location displacements and covariances from a sequence of gravity-aligned inertial data. To avoid the impacts from initial orientation, the inertial data are transformed into a local gravity-aligned frame. The learned displacements and covariances are then incorporated into an extended Kalman filter as observation states that estimate full-states of orientation, velocity, location, and IMU bias. In a 3-7 minute human

motion scenario, the localization error of TLIO is within 3 meters 90% of the time.

Another category of deep learning models aims to leverage learned velocity to correct accelerations, as illustrated in Figure 5(b). A prominent example is RIDI [35], which trains a deep neural network to predict velocity vectors from inertial data, which are then used to correct linear accelerations by subtracting gravity, aligning with the constraints of learned velocities. The corrected linear accelerations are then doubly integrated to estimate positions. To enhance the accuracy of inertial accelerations, RIDI leverages human walking speed as a prior, which compensates for the drifts in inertial positioning, effectively constraining them to a lower level. RoNIN [46] improves upon RIDI by transforming inertial measurements and learned velocity vectors into a heading-agnostic coordinate frame and introducing several novel velocity losses. To minimize the impact of orientation estimation, RoNIN employs device orientation to transform inertial data into a frame with its Z-axis aligned with gravity. However, a limitation of RoNIN is its reliance on orientation estimation. NILoc [56] is an intriguing trial based on RoNIN, which tackles the neural inertial localization problem, aiming to infer global location from inertial motion history only. This work recognizes that human motion patterns are unique in different locations, which can be utilized as a "fingerprint" to determine the location, similar to WiFi or magnetic-field fingerprinting. NILoc first calculates a sequence of velocity from inertial data and then employs a Transformer-based DNN framework [72] to transform the velocity sequence into location. However, one fundamental limitation of NILoc is that in some areas, such as open spaces, symmetrical or repetitive places, there may not be a unique motion pattern.

An alternative approach involves incorporating learned velocity into the updating process of a Kalman filter (KF), as shown in Figure 5 (c). [36] uses a ConvNet to infer current speed from IMU sequences and incorporates this speed into the Kalman filter as a velocity observation to constrain the drifts of SINS-based inertial positioning. This approach is similar to the zero-velocity update (ZUPT) method, which detects and uses zero-velocity in KF as observations, but instead uses full speeds as observations in KF. Incorporating learned velocity allows the KF to handle more complex human motion. A similar trial is [49], that is based on a DNN that infers walking velocity in the body frame and combines it with an extended KF. In addition to the learned velocity, [49] produces a noise parameter for KF to dynamically update parameters, rather than setting a fixed noise parameter.

Inertial positioning heavily relies on accurately estimating the device's attitude. Several methods aim to improve orientation estimation to enhance the performance of deep learning based inertial odometry. RIDI, RoNIN, and TLIO still depend on device orientation to rotate inertial data into a suitable frame. To address this problem, IDOL [51] proposes a two-stage process that first learns orientation from data and then rotates inertial data into the appropriate frame, followed by learning the position. [58] estimates orientation using magnetic data and combines it with learned odometry to reduce positioning drifts while minimizing reliance on device orientation.
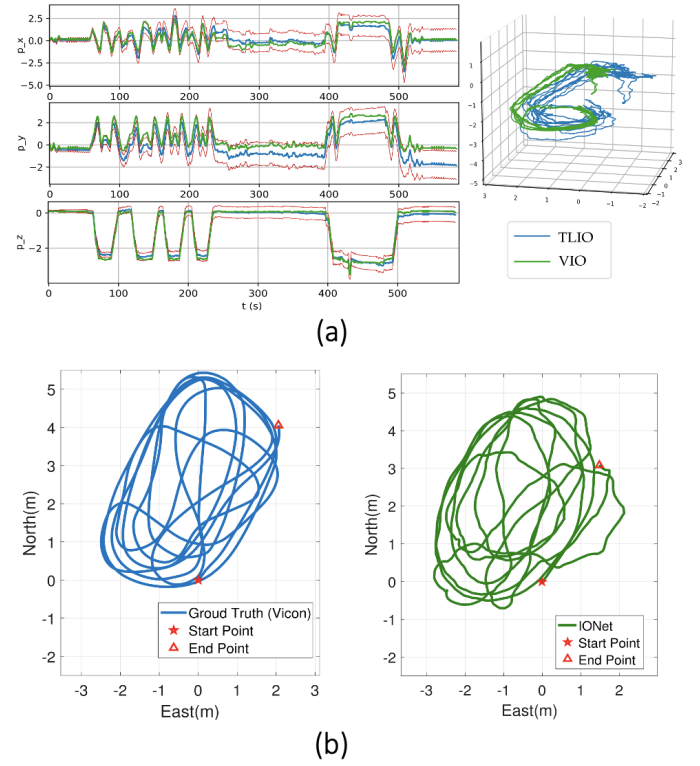


Fig. 6: Sample results of deep learning based inertial positioning from (a) VR device for pedestrian tracking (reprint from TLIO [44]) (b) smartphone for trolly tracking (reprint from IONet [34])

Figure 6 showcases several examples of deep learning based inertial positioning results.

## V. LEARNING TO CORRECT PEDESTRIAN INERTIAL POSITIONING

The previous subsection addressed the general application of deep learning in correcting inertial positioning drifts. This subsection focuses on the specific use of deep learning to address particular aspects of pedestrian navigation algorithms, namely Pedestrian Dead Reckoning (PDR) and Zero-Velocity Update (ZUPT).

### A. Learning to correct PDR

Pedestrian dead reckoning (PDR) error drifts often stem from inaccurate stride and heading estimates. To address these issues, researchers have incorporated deep learning techniques into the process of step detection, dynamic step length estimation, and walking heading estimation.

To estimate walking stride more robustly, researchers have sought to solve it in a data-driven way. One such method is SmartStep [73], a deep learning-based step detection framework that achieves 99% accuracy in step detection tasks across various motion modes. Compared to peak/valley detection-based methods, data-driven methods do not require IMUs to be fixed in position, specific motion modes, or pre-calibration and threshold setting. Another approach involves using LSTM to

regress walking stride from raw inertial data [41]. This method has demonstrated effectiveness in various human motions, such as walking, running, jogging, and random movements. Additionally, StepNet [48] learns to estimate step length dynamically, i.e., the change in distance, which achieves an impressive performance with only a 2.1%-3.2% error rate when compared to traditional static step length estimation. The attachment mode of the device, such as in hand or in pocket, can also influence walking stride estimation. To address this problem, Bo et al. [62] employed domain adaptation [70] to extract domain-invariant features for stride estimation, which enhanced the performance in new domains, such as holding, calling, pocket, and swinging.

Accurate heading estimation is crucial for updating position in the right direction in PDR. To achieve more accurate and robust heading estimation, Wang et al. [42] utilize a Spatial Transformer Network [74] and LSTM to learn heading direction from the inertial sensor attached to an unconstrained device. However, one problem that arises is the misalignment between the device heading and pedestrian heading, making it difficult to estimate the real walking heading based on sensor data. To address this misalignment issue, [75] introduces a deep neural network to estimate walking direction in the sensor's frame. They derive a geometric model to convert walking direction from the sensor's frame into a reference frame (i.e., north and east coordinates) by exploiting acceleration and magnetic data. This geometric model is combined with a learning framework to produce heading estimates. When tested on unseen data, this work reports a median heading error of $10°$. PDRNet [52] follows the process of a traditional PDR algorithm but replaces the step length and heading estimation modules with deep neural networks. Their experiments indicate that learning step length and heading together outperforms regressing them separately.

### B. Learning to correct ZUPT

In pedestrian inertial navigation systems (INS) based on zero-velocity update (ZUPT), the zero-velocity phase is utilized to correct inertial positioning errors through Kalman filtering. Therefore, the accuracy of zero-velocity detection is crucial in determining when to update the system states. However, traditional threshold-based zero-velocity detection is complicated by the mixed variety of motions experienced by humans, making it challenging to set a reliable threshold when the user is still.

To address this issue, researchers have explored data-driven approaches that utilize the powerful feature extraction and classification capabilities of deep learning to classify whether the user is in the ZUPT phase. For instance, [37] proposes a six-layer long short-term memory (LSTM) network to detect zero-velocity. The LSTM inputs a sequence of IMU data, typically 100 consecutive data points, and outputs the probability of whether the user is still or in motion at the current timestep. The results from the LSTM-based zero-velocity detection are then fed into a ZUPT-based INS. The proposed approach achieves a reduction in localization error by over 34% compared to fixed threshold-based ZVDs and was shown to be more robust during a mixed variety of motions, such as walking, running, and climbing stairs. Similarly, [43] designs an adaptive ZUPT using convolutional neural networks (ConvNet) to classify ZVDs based on IMU sequences. Deep learning approaches, such as LSTM and ConvNet, have demonstrated excellent performance in extracting robust and useful features for zero-velocity identification, irrespective of different users, motion modes, and attachment places.

## VI. LEARNING TO CORRECT INERTIAL POSITIONING ON VEHICLES, UAV AND ROBOTIC PLATFORMS

As previously mentioned, deep learning methods have shown great potential in addressing the challenges of pedestrian inertial navigation. However, these techniques can also be applied to other platforms, such as vehicles, UAVs, robots, and more.

These platforms share similarities with pedestrians, such as the ability to infer movement velocity from inertial data. This is because inertial data contains vibration information that reflects the fundamental frequency proportional to the vehicle speed. Building on the success of IONet [34], [39] proposes AbolDeepIO, an improved triple-channel LSTM network that predicts polar vectors for drone localization from inertial data sequences. AbolDeepIO has been evaluated on a public drone dataset and has shown competitive performance compared to traditional visual-inertial odometry methods like VINS-mono.

When deploying deep learning-based inertial navigation on real-world devices, prediction accuracy and model efficiency must be considered. To address this, TinyOdom [59] aims to deploy neural inertial odometry models on resource-constrained devices. It proposes a lightweight model based on temporal convolutional networks (TCN) [76] to learn position displacement and optimizes the model through neural architecture search (NAS) [77] to reduce model size between 31 and 134 times. TinyOdom was extensively evaluated on tracking pedestrians, animals, aerial, and underwater vehicles. Within 60 seconds, its localization error is between 2.5 and 12 meters.

Learning-based inertial odometry has also been extended to legged robots by [53]. The learned location displacement is combined with kinematic motion models to estimate robot system states at high frequencies (400 Hz). In this work, the robot successfully navigated a field experiment, where a legged robot walked around for 20 minutes in a mine with poor illumination and visual feature tracking failures.

In the realm of inertial positioning for vehicles, researchers have proposed various methods to mitigate error drifts and improve accuracy. One such method is presented in [47], where error covariances are learned from inertial data and incorporated into Kalman filtering for updating system states. This approach has been shown to improve inertial positioning performance. Similar to ZUPT-based pedestrian positioning, zero-velocity-update (ZUPT) can also be used for car-equipped inertial navigation systems. The zero-velocity phase provides valuable context information to correct system error drifts via Kalman filtering. OdoNet, presented in [63], is an example of a system that learns and utilizes car speed along with a zero-velocity detector to reduce error drifts in car-equipped IMU

TABLE III: A summary of existing methods on deep learning based sensor fusion.

| name | year | sensor | model | learning | target |
|---|---|---|---|---|---|
| VINet [78] | 2017 | MC+I | ConvNet, LSTM | SL | formulating VIO as a sequential learning problem |
| VIOLearner [79] | 2018 | MC+I | ConvNet | UL | VIO with online correction module |
| Chen et al. [80] | 2019 | MC+I | ConvNet, LSTM, Attention | SL | feature selection for deep VIO |
| DeepVIO [81] | 2019 | SC+I | ConvNet, LSTM | UL | learning VIO from stereo images and IMU |
| DeepTIO [82] | 2020 | T+I | ConvNet, LSTM, Attention | SL | learning pose from thermal and inertial data |
| MilliEgo [83] | 2020 | MR+I | ConvNet, LSTM, Attention | SL | learning pose from mmWare radar and inertial data |
| UnVIO [84] | 2021 | MC+I | ConvNet, LSTM, Attention | UL | unsupervised learning of VIO |
| DynaNet [85] | 2021 | MC+I | ConvNet, LSTM | SL | combining DNN with Kalman filtering |
| SelfVIO [86] | 2022 | MC+I | ConvNet, LSTM, Attention | UL | unsupervised VIO with GAN-based depth generator |
| Tu et al. [87] | 2022 | L+I | ConvNet, LSTM, Attention | UL | unsupervised learning of LIDAR-inertial odometry |

- *Year* indicates the publication year of each work.
- *Sensor* indicates the sensors involved in each work. I, MC, SC, T, MR, L, A represent inertial sensor, monocular camera, stereo camera, thermal camera, millimeter wave radar, LIDAR and airflow sensor respectively.
- *Learning* indicates how to train neural networks. SL and UL represent Supervised Learning and Unsupervised Learning.

systems. Deep learning techniques have also been explored for detecting zero-velocity phases in vehicle navigation. For example, [40] proposes a deep learning-based method for detecting zero-velocity phases in vehicle navigation. In another study, [54] derives a model with motion terms that are relevant only to the IMU data sequence. This model provides theoretical guidance for learning models to infer useful terms and has been evaluated on a drone dataset, where it outperformed TLIO and other learning methods.

Overall, these studies demonstrate the potential of deep learning-based methods in improving inertial navigation for various platforms, including pedestrians, vehicles, drones, and robots. By leveraging the rich information contained within IMU data, deep learning models can effectively mitigate error drifts and improve the accuracy of inertial positioning systems. Furthermore, by optimizing the model efficiency and considering deployment on resource-constrained devices, these techniques can be applied in real-world scenarios.

## VII. DEEP LEARNING BASED SENSOR FUSION

Integrating inertial sensors with other sensors as a multi-sensor navigation system has been an area of research for several decades. Nowadays, platforms such as robots, vehicles, and VR/AR devices are equipped with cameras, IMUs, and LIDAR sensors. Hence, it is natural to consider introducing multimodal learning techniques [88] and designing learning models capable of fusing multimodal information to construct a mapping function from sensor data to pose.

Visual-inertial odometry (VIO) has garnered attention as a means of integrating low-cost, complementary camera and IMU sensors that are widely deployed. Monocular vision can capture the appearance and geometry of a scene, but cannot recover the scale metric. IMU provides metric scale and improves motion tracking in featureless areas, complex lighting conditions, and motion blur. However, a pure inertial solution can only last for a short period. Therefore, an effective fusion of these two complementary sensors is necessary for accurate pose estimation.

Traditional VIO methods integrate visual and inertial information based on filtering [5], [89], fixed-lag smoothing [90], or full smoothing [91]. Recently, deep learning-based VIO
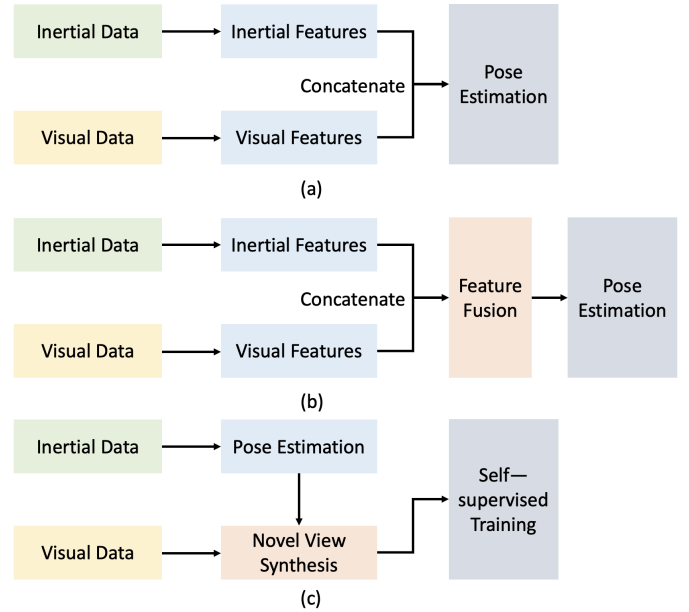


Fig. 7: An overview of existing methods on deep learning based sensor fusion

models have emerged, directly constructing a mapping function from images and IMU to pose in a data-driven manner. VINet [78] is an end-to-end deep VIO model consisting of a ConvNet-based visual encoder to extract visual features from two images and an LSTM-based inertial encoder to extract inertial features from a sequence of inertial data between the two images. As shown in Figure 7 (a), the visual and inertial features are concatenated together as one tensor, followed by an LSTM and fully-connected layer that finally maps features into a 6-dimensional pose. VINet is trained on public driving datasets such as the KITTI dataset [92] and a public drone dataset such as the EuroC dataset [93]. The learned VIO model is generally more robust to sensor noises compared to traditional VIO methods, although its model performance still cannot compete with state-of-the-art VIO methods.

To effectively integrate visual and inertial information, [80] proposes a selective sensor fusion mechanism that learns to choose important features conditioned on sensor observations,

as demonstrated in Figure 7 (b). Specifically, this work proposes two types of fusion: soft fusion, which is based on an attention mechanism and generates a soft mask to reweight features based on their importance, and hard fusion, which is based on Gumbel Soft-max and generates a hard mask consisting of either 1 or 0 to either propagate or ignore a feature. Experimental evaluation on the KITTI dataset demonstrates that compared with directly concatenating features [78], selective fusion enhances the performance of deep VIO by 5%-10%. An interesting observation is that the number of useful features is relevant to the amount of linear/rotational velocity, with inertial features contributing more to rotation rate (e.g., turning), while more visual features are used to increase linear velocity.

Both [78] and [80] are trained in a supervised learning manner using datasets with high-precision ground-truth poses as training labels. However, obtaining high-precision poses can be difficult or costly in certain cases. Consequently, self-supervised learning-based VIOs, which do not require pose labels, have attracted attention. Self-supervised VIOs leverage the multi-view geometry relation of consecutive images, such as novel view synthesis, as a supervision signal [79], [81], [84], [86]. The task of novel view synthesis involves transforming a source image into a target view and comparing the differences between the synthesized target images and real target images as loss. In VIOLearner [79] and DeepVIO [81], as shown in Figure 7 (c), the pose transformation is generated from an inertial data sequence and used in the novel view synthesis process. In UnVIO [84] and SelfVIO [86], inertial data is integrated with visual data via an attention module applied to the concatenated visual and inertial features extracted from the images and IMU sequence. They show that incorporating inertial data with visual data improves the accuracy of pose estimation, particularly rotation estimation.

The use of learning-based sensor fusion extends beyond visual-inertial odometry (VIO) to include other sensor modalities such as Lidar-inertial odometry (LIO), thermal-inertial odometry, and radar-inertial odometry [82], [83], [87]. Deep-TIO [82] and MilliEgo [83] employ attention-based selective fusion mechanisms, similar to soft fusion [80], to reweight and fuse features from inertial and visual data, resulting in improved pose accuracy. In addition, unsupervised learning-based LIDAR-inertial odometry [87] generates motion transformation from IMU sequence and uses it for LIDAR novel view synthesis to facilitate self-supervised learning of ego-motion, similar to VIOLearner [79]. In all these cases, the inclusion of IMU data in deep neural networks enhances pose estimation accuracy and robustness.

## VIII. DEEP LEARNING BASED HUMAN MOTION ANALYSIS AND ACTIVITY RECOGNITION

Inertial sensors have diverse applications beyond positioning, such as motion tracking, activity recognition, and more. Although these tasks are not the primary focus of this survey, this section provides a brief yet comprehensive overview of how deep learning is utilized in these domains.

### A. Human Motion Analysis

Data-driven approaches are utilized to reconstruct human pose and motion using either a single IMU or multiple IMUs attached to the body. These models primarily focus on analyzing human motion rather than localizing users, which differentiates them from inertial positioning. Several studies have applied machine learning to gait and pose analysis, such as knee angle estimation for human walking using supervised support vector regression in [94] and probabilistic parameter learning for human gesture recognition in [95] through handcrafted motion features extracted from inertial data. In addition, machine learning methods, such as multi-layer perceptrons (MLPs), have been utilized in IMU data to learn sensor displacement for human motion reconstruction in [96]–[98].

Recently, deep learning has shown promising performance in human pose reconstruction. For example, [99] proposed Deep Inertial Poser, a recurrent neural network (RNN)-based framework that can reconstruct full-body pose from six IMUs attached to the user's body. TransPose [100], another RNN-based framework, enables real-time human pose estimation using six body-attached IMUs. Furthermore, [101] combines a neural kinematics estimator with a physics-aware motion optimizer to improve the accuracy of human motion tracking.

### B. Human Activity recognition (HAR)

Deep learning can be utilized to exploit inertial information from body-worn IMUs for human activity recognition. For instance, [102] published a popular public dataset of human activity recognition and successfully classified current activity among six classes, including walking, standing still, sitting, walking downstairs, walking upstairs, and laying down, using support vector machines (SVM). In addition, [103] presents an LSTM-based HAR model that inputted a sequence of inertial data and outputted class probability. Moreover, [104] introduces a ConvNet-based HAR model that achieved a classification accuracy of 97%, outperforming an accuracy of 96% from SVM-based HAR models. To reduce onboard computational requirements, [105] presents a learning framework that exploited both features automatically extracted by DNN and hand-crafted features to achieve accurate and real-time human activity recognition on low-end devices.

Learning from inertial data can also benefit sports and health applications. For instance, [106] shows that deep learning is effective in detecting Parkinson's disease by assessing the patient's daily activity through the analysis of inertial information from wearable sensors. Additionally, [107] provides instructions for athletes' sports training based on sensor data and activity information.

## IX. CONCLUSIONS AND DISCUSSIONS

In recent years, there has been a growing interest in using deep learning to address the problem of inertial positioning. This article provides a comprehensive review of the area of deep learning-based inertial positioning. The rapid advances in this field have already provided promising solutions to address problems such as inertial sensor calibration, the compensation

of error drifts in inertial positioning, and multimodal sensor fusion. This section concludes and discusses the benefits that deep learning can bring to inertial navigation research, analyzes the challenges that existing research faces, and highlights the future opportunities of this evolving field.

### A. Benefits

Unlike traditional geometric or physical inertial positioning models, the integration of deep learning into inertial positioning has led to the development of a range of alternative solutions to address the issue of positioning error drifts. The corresponding benefits can be summarized as follows:

*1) Learn to approximate complex and varying function:* The deep neural network has proven to be a powerful and versatile nonlinear function that can approximate the complex and variable factors involved in inertial positioning, which are difficult to model manually. For example, when calibrating sensors, the corrupt noises that exist in inertial measurements can be modeled and eliminated in a data-driven way by training on a large dataset using a DNN. Deep learning can also directly generate absolute velocity and position displacement from data, without the need for IMU integration, thus reducing positioning drifts. In pedestrian dead reckoning (PDR), deep learning can estimate step length based on data, rather than empirical equations, and implicitly remove the effects of different users. These works demonstrate that using a large dataset to build a data-driven model can produce more accurate motion estimates, as well as reduce and constrain the rapid error drifts of inertial navigation systems.

*2) Learn to estimate parameters:* Automatic identification of parameters through data-driven models contributes to paving the way for next-generation intelligent navigation systems that can actively exploit input data and evolve over time without human intervention. In classical inertial navigation mechanisms, certain parameters or modules need to be manually set and tuned before use. For instance, experts with experience need to settle parameters in Kalman filtering, such as observation noise, covariance, and process noise. Deep learning has proven effective in automatically producing suitable parameters for Kalman filtering based on input data [42], [47], [85]. In sensor calibration, reinforcement learning algorithms are used to discover optimal parameters for inertial calibration algorithms [27]. In ZUPT-based pedestrian inertial positioning, deep learning is a viable solution for classifying zero-velocity phases and determining when to update system states.

*3) Learn to self-adapt in new domains:* Unforeseen or ever-changing issues in new application domains, such as changes in motion mode, carrier, and sensor noise, can significantly impact the performance of inertial systems. Learning models offer opportunities for inertial systems to adapt to new changes and overcome these influential factors implicitly by discovering and exploiting the differences in data distributions between domains. For instance, [38] leverages transfer learning to allow INS to extract domain-invariant features from data, maintaining localization accuracy when sensor attachment is changed. The introduction of self-supervised learning enables navigation systems to learn from data without high-precision pose as training labels, allowing unlabelled inertial data to be effectively used for model performance improvement. In visual-inertial odometry, [79], [81], [84] introduce novel view synthesis as a supervision signal to train deep VIO in a self-supervised learning way. This self-adaptation ability is promising for mobile agents to continuously improve their localization performance in new application scenes.

### B. Challenges and Opportunities

Despite the impressive and promising results that deep learning has already offered in inertial positioning, there are still challenges in existing methods when they are applied and deployed in real-world scenarios. To overcome these limitations, several opportunities and potential research directions are discussed below.

*1) Generalization and Self-learning:* The generalization problem is a major concern for deep learning-based methods because these models are trained on one domain (i.e., training set) but need to be tested on other domains (i.e., testing set). The possible differences in data between domains can lead to a degradation of prediction performance. Although deep learning-based inertial navigation models have reported impressive results on the author's own datasets, these works have not been evaluated in comprehensive experiments during long-term operation and across various devices, users, and application scenes. Thus, it is challenging to determine the real performance of these models in open environments. To address the generalization problem, new learning techniques such as transfer learning [108], lifelong learning, and contrastive learning [109] can be introduced into inertial positioning systems, which is a promising direction. For instance, in the future, by exploiting information from physical/geometric rules or other sensors (e.g., GNSS, camera), the learning-based inertial positioning model can be self-supervisedly trained and enable mobile agents to learn from data in a lifelong manner.

*2) Black-box and Explainability:* Deep neural networks have been criticized as being a 'black-box' model due to their lack of explainability and interpretability. As these models are often used to support real-world tasks, it is crucial to investigate what is learned inside deep nets before deploying them to ensure their safety and reliability. Despite the good results shown by deep learning models in estimating important terms such as location displacement, sensor measurement errors, and filtering parameters, these terms lack concrete mathematical models, unlike traditional inertial navigation. To determine whether these terms are trustworthy, uncertainties should be estimated in conjunction with the inertial positioning method [71] and used as indicators for users or systems to understand the extent to which model predictions can be trusted. In future research, it is important to reveal the governing mathematical or physical models behind the learned inertial positioning neural model and identify which parts of inertial positioning can be learned by deep nets. Introducing Bayesian deep learning into inertial positioning is also a promising direction that could offer interpretability for model predictions [110].

*3) Efficiency and Real-world Deployment:* When deploying deep positioning models on user devices, it is crucial to consider the consumption of computation, storage, and energy in system design, in addition to prediction accuracy. Compared to classical inertial navigation algorithms, DNN-based inertial positioning models have a relatively large computational and memory burden, as they contain millions of neural parameters that require GPUs for parallel training and testing. Therefore, online inference of learning models, especially on low-end devices such as IoT consoles, VR/AR devices, and miniature drones, requires lightweight, efficient, and effective models. To achieve this goal, neural model compression techniques, such as knowledge distillation [111], should be introduced to discover the optimal neural structure that balances prediction accuracy and model size. [45] and [63] have conducted initial trials on minimizing the model size of inertial odometry. Moreover, safety and reliability are also crucial factors to consider. In the future, it is worth exploring the optimal structure of learning-based inertial positioning models, considering model performance, parameter size, latency, safety, and reliability for real-world deployment.

*4) Data Collection and Benchmark:* As a data-driven approach, the performance of deep learning models depends on the quality of the data, such as the size of the dataset, data diversity, and the differences between the training and testing sets. Under ideal conditions, deep learning-based inertial positioning models should be trained on diverse data across different users, platforms, motion dynamics, and inertial sensors to enhance their generalization in testing domains. However, collecting such data in diverse domains can be costly and time-consuming. Additionally, obtaining high-precision ground-truth poses as training and evaluation labels can be challenging in some cases. Previous research has used different training/evaluation data, model hyperparameters (e.g., learning rate, batch size, layer dimension), and evaluation metrics, making it difficult to compare these methods fairly. In visual navigation tasks, such as visual odometry/SLAM, the KITTI dataset [92] is commonly used as a benchmark to train and evaluate learning-based VO models. However, although published datasets for inertial navigation exist [45], [46], there is still a lack of a common benchmark that is adopted and recognized by mainstream methods in inertial positioning. In the future, a widely adopted dataset and benchmark, covering a variety of application scenarios, will greatly benefit and foster research in data-driven inertial positioning.

*5) Failure Cases and Physical Constraints:* Deep learning has demonstrated its capability in reducing the drifts of inertial positioning and contributing to various aspects of inertial navigation systems, as discussed in Section IV. However, DNN models are not always reliable and may occasionally produce large and abrupt prediction errors. Unlike traditional inertial navigation algorithms that are based on concrete physical and mathematical rules, DNN predictions lack constraints, and the failure cases must be considered in real-world applications with safety concerns. To enhance the robustness of DNN predictions, possible solutions include imposing physical constraints on DNN models or combining deep learning with physical models as hybrid inertial positioning models. By doing so, the benefits from both learning and physics-based positioning models can be leveraged.

*6) New deep learning methods:* Machine/deep learning is one of the fastest growing areas of AI, and its advances have influenced numerous fields such as computer vision, robotics, natural language processing, and signal processing. There are significant opportunities for applying deep learning techniques to inertial navigation and analyzing their effectiveness and theoretical underpinnings. In the future, new model structures such as transformer [72], diffusion models [112], and generative models [69], and new learning methods such as transfer learning, reinforcement learning, contrastive learning [109], unsupervised learning, and meta-learning [113], all hold promise for enhancing inertial positioning systems. Furthermore, advances in other domains such as neural rendering [114] and voice synthesis [115] may provide valuable insights into developing more effective inertial positioning systems. Therefore, incorporating these rapidly-evolving deep learning methods into inertial navigation will be a significant area of research in the future.

## REFERENCES

[1] P. G. Savage, "Strapdown Inertial Navigation Integration Algorithm Design Part 1: Attitude Algorithms," *Journal of Guidance, Control, and Dynamics*, vol. 21, no. 1, pp. 19–28, 1998.

[2] P. G. Savage, "Strapdown Inertial Navigation Integration Algorithm Design Part 2: Velocity and Position Algorithms," *Journal of Guidance, Control, and Dynamics*, vol. 21, no. 1, pp. 19–28, 1998.

[3] R. Harle, "A survey of indoor inertial positioning systems for pedestrians," *IEEE Communications Surveys & Tutorials*, vol. 15, no. 3, pp. 1281–1293, 2013.

[4] I. Skog, P. Händel, J.-O. Nilsson, and J. Rantakokko, "Zero-Velocity Detection — an Algorithm Evaluation.," *IEEE transactions on biomedical engineering*, vol. 57, no. 11, pp. 2657–2666, 2010.

[5] M. Li and A. I. Mourikis, "High-precision, Consistent EKF-based Visual-Inertial Odometry," *The International Journal of Robotics Research*, vol. 32, no. 6, pp. 690–711, 2013.

[6] T. Qin, P. Li, and S. Shen, "VINS-Mono: A Robust and Versatile Monocular Visual-Inertial State Estimator," *IEEE Transactions on Robotics*, vol. 34, no. 4, pp. 1004–1020, 2018.

[7] W. Xu, Y. Cai, D. He, J. Lin, and F. Zhang, "Fast-lio2: Fast direct lidar-inertial odometry," *IEEE Transactions on Robotics*, 2022.

[8] Y. Bengio, I. Goodfellow, and A. Courville, *Deep learning*, vol. 1. MIT press Cambridge, MA, USA, 2017.

[9] Z.-Q. Zhao, P. Zheng, S.-t. Xu, and X. Wu, "Object detection with deep learning: A review," *IEEE transactions on neural networks and learning systems*, vol. 30, no. 11, pp. 3212–3232, 2019.

[10] S. Hao, Y. Zhou, and Y. Guo, "A brief survey on semantic segmentation with deep learning," *Neurocomputing*, vol. 406, pp. 302–321, 2020.

[11] N. Sünderhauf, O. Brock, W. Scheirer, R. Hadsell, D. Fox, J. Leitner, B. Upcroft, P. Abbeel, W. Burgard, M. Milford, *et al.*, "The limits and potentials of deep learning for robotics," *The International journal of robotics research*, vol. 37, no. 4-5, pp. 405–420, 2018.

[12] Y. Li, R. Chen, X. Niu, Y. Zhuang, Z. Gao, X. Hu, and N. El-Sheimy, "Inertial sensing meets machine learning: Opportunity or challenge?," *IEEE Transactions on Intelligent Transportation Systems*, 2021.

[13] P. S. Farahsari, A. Farahzadi, J. Rezazadeh, and A. Bagheri, "A survey on indoor positioning systems for iot-based applications," *IEEE Internet of Things Journal*, vol. 9, no. 10, pp. 7680–7699, 2022.

[14] L. E. Díez, A. Bahillo, J. Otegui, and T. Otim, "Step length estimation methods based on inertial sensors: A review," *IEEE Sensors Journal*, vol. 18, no. 17, pp. 6908–6926, 2018.

[15] Y. Wu, H.-B. Zhu, Q.-X. Du, and S.-M. Tang, "A survey of the research status of pedestrian dead reckoning systems based on inertial sensors," *International Journal of Automation and Computing*, vol. 16, no. 1, pp. 65–83, 2019.

[16] X. Ru, N. Gu, H. Shang, and H. Zhang, "Mems inertial sensor calibration technology: Current status and future trends," *Micromachines*, vol. 13, no. 6, p. 879, 2022.

[17] N. Naser, El-Sheimy; Haiying, Hou; Xiaojii, "Analysis and Modeling of Inertial Sensors Using Allan Variance," *IEEE Transactions on Instrumentation and Measurement*, vol. 57, no. JANUARY, pp. 684–694, 2008.

[18] H. Weinberg, "Using the adxl202 in pedometer and personal navigation applications," *Analog Devices AN-602 application note*, vol. 2, no. 2, pp. 1–6, 2002.

[19] L. Fang, P. J. Antsaklis, L. A. Montestruque, M. B. McMickell, M. Lemmon, Y. Sun, H. Fang, I. Koutroulis, M. Haenggi, M. Xie, *et al.*, "Design of a wireless assisted pedestrian dead reckoning system-the navmote experience," *IEEE transactions on Instrumentation and Measurement*, vol. 54, no. 6, pp. 2342–2358, 2005.

[20] P. Goyal, V. J. Ribeiro, H. Saran, and A. Kumar, "Strap-down pedestrian dead-reckoning system," in *2011 international conference on indoor positioning and indoor navigation*, pp. 1–7, IEEE, 2011.

[21] B. Huang, G. Qi, X. Yang, L. Zhao, and H. Zou, "Exploiting cyclic features of walking for pedestrian dead reckoning with unconstrained smartphones," in *Proceedings of the 2016 ACM International Joint Conference on Pervasive and Ubiquitous Computing*, pp. 374–385, 2016.

[22] D. Feng, C. Wang, C. He, Y. Zhuang, and X.-G. Xia, "Kalman-filter-based integration of imu and uwb for high-accuracy indoor positioning and navigation," *IEEE Internet of Things Journal*, vol. 7, no. 4, pp. 3133–3146, 2020.

[23] S. Yang, J. Liu, X. Gong, G. Huang, and Y. Bai, "A robust heading estimation solution for smartphone multisensor-integrated indoor positioning," *IEEE Internet of Things Journal*, vol. 8, no. 23, pp. 17186–17198, 2021.

[24] C. Xiyuan, "Modeling random gyro drift by time series neural networks and by traditional method," in *International Conference on Neural Networks and Signal Processing, 2003. Proceedings of the 2003*, vol. 1, pp. 810–813, IEEE, 2003.

[25] H. Chen, P. Aggarwal, T. M. Taha, and V. P. Chodavarapu, "Improving inertial sensor by reducing errors using deep learning methodology," in *NAECON 2018-IEEE National Aerospace and Electronics Conference*, pp. 197–202, IEEE, 2018.

[26] M. A. Esfahani, H. Wang, K. Wu, and S. Yuan, "Orinet: Robust 3-d orientation estimation with a single particular imu," *IEEE Robotics and Automation Letters*, vol. 5, no. 2, pp. 399–406, 2019.

[27] F. Nobre and C. Heckman, "Learning to calibrate: Reinforcement learning for guided calibration of visual–inertial rigs," *The International Journal of Robotics Research*, vol. 38, no. 12-13, pp. 1388–1402, 2019.

[28] M. Brossard, S. Bonnabel, and A. Barrau, "Denoising imu gyroscopes with deep learning for open-loop attitude estimation," *IEEE Robotics and Automation Letters*, vol. 5, no. 3, pp. 4796–4803, 2020.

[29] X. Zhao, C. Deng, X. Kong, J. Xu, and Y. Liu, "Learning to compensate for the drift and error of gyroscope in vehicle localization," in *2020 IEEE Intelligent Vehicles Symposium (IV)*, pp. 852–857, IEEE, 2020.

[30] F. Huang, Z. Wang, L. Xing, and C. Gao, "A mems imu gyroscope calibration method based on deep learning," *IEEE Transactions on Instrumentation and Measurement*, vol. 71, pp. 1–9, 2022.

[31] R. Li, C. Fu, W. Yi, and X. Yi, "Calib-net: Calibrating the low-cost imu via deep convolutional neural network," *Frontiers in Robotics and AI*, vol. 8, p. 772583, 2022.

[32] S.-i. Amari, "Backpropagation and stochastic gradient descent method," *Neurocomputing*, vol. 5, no. 4-5, pp. 185–196, 1993.

[33] A. K. Jain, J. Mao, and K. M. Mohiuddin, "Artificial neural networks: A tutorial," *Computer*, vol. 29, no. 3, pp. 31–44, 1996.

[34] C. Chen, X. Lu, A. Markham, and N. Trigoni, "Ionet: Learning to cure the curse of drift in inertial odometry," in *The Conference on Artificial Intelligence (AAAI)*, 2018.

[35] H. Yan, Q. Shan, and Y. Furukawa, "Ridi: Robust imu double integration," in *Proceedings of the European Conference on Computer Vision (ECCV)*, pp. 621–636, 2018.

[36] S. Cortés, A. Solin, and J. Kannala, "Deep learning based speed estimation for constraining strapdown inertial navigation on smartphones," in *2018 IEEE 28th International Workshop on Machine Learning for Signal Processing (MLSP)*, pp. 1–6, IEEE, 2018.

[37] B. Wagstaff and J. Kelly, "Lstm-based zero-velocity detection for robust inertial navigation," in *2018 International Conference on Indoor Positioning and Indoor Navigation (IPIN)*, pp. 1–8, IEEE, 2018.

[38] C. Chen, Y. Miao, C. X. Lu, L. Xie, P. Blunsom, A. Markham, and N. Trigoni, "Motiontransformer: Transferring neural inertial tracking between domains," in *The Conference on Artificial Intelligence (AAAI)*, vol. 33, 2019.

[39] M. A. Esfahani, H. Wang, K. Wu, and S. Yuan, "Aboldeepio: A novel deep inertial odometry network for autonomous vehicles," *IEEE Transactions on Intelligent Transportation Systems*, 2019.

[40] M. Brossard, A. Barrau, and S. Bonnabel, "Rins-w: Robust inertial navigation system on wheels," *The IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, 2019.

[41] T. Feigl, S. Kram, P. Woller, R. H. Siddiqui, M. Philippsen, and C. Mutschler, "A bidirectional lstm for estimating dynamic human velocities from a single imu," in *2019 International Conference on Indoor Positioning and Indoor Navigation (IPIN)*, pp. 1–8, IEEE, 2019.

[42] Q. Wang, H. Luo, L. Ye, A. Men, F. Zhao, Y. Huang, and C. Ou, "Pedestrian heading estimation based on spatial transformer networks and hierarchical lstm," *IEEE Access*, vol. 7, pp. 162309–162322, 2019.

[43] X. Yu, B. Liu, X. Lan, Z. Xiao, S. Lin, B. Yan, and L. Zhou, "Azupt: Adaptive zero velocity update based on neural networks for pedestrian tracking," in *2019 IEEE Global Communications Conference (GLOBECOM)*, pp. 1–6, IEEE, 2019.

[44] W. Liu, D. Caruso, E. Ilg, J. Dong, A. I. Mourikis, K. Daniilidis, V. Kumar, and J. Engel, "Tlio: Tight learned inertial odometry," *IEEE Robotics and Automation Letters*, vol. 5, no. 4, pp. 5653–5660, 2020.

[45] C. Chen, P. Zhao, C. X. Lu, W. Wang, A. Markham, and N. Trigoni, "Deep-learning-based pedestrian inertial navigation: Methods, data set, and on-device inference," *IEEE Internet of Things Journal*, vol. 7, no. 5, pp. 4431–4441, 2020.

[46] S. Herath, H. Yan, and Y. Furukawa, "Ronin: Robust neural inertial navigation in the wild: Benchmark, evaluations, & new methods," in *2020 IEEE International Conference on Robotics and Automation (ICRA)*, pp. 3146–3152, IEEE, 2020.

[47] M. Brossard, A. Barrau, and S. Bonnabel, "Ai-imu dead-reckoning," *IEEE Transactions on Intelligent Vehicles*, vol. 5, no. 4, pp. 585–595, 2020.

[48] I. Klein and O. Asraf, "Stepnet—deep learning approaches for step length estimation," *IEEE Access*, vol. 8, pp. 85706–85713, 2020.

[49] Y. Wang, H. Cheng, and M. Q.-H. Meng, "Pedestrian motion tracking by using inertial sensors on the smartphone," in *2020 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, pp. 4426–4431, IEEE, 2020.

[50] X. Teng, P. Xu, D. Guo, Y. Guo, R. Hu, H. Chai, and D. Chuxing, "Arpdr: An accurate and robust pedestrian dead reckoning system for indoor localization on handheld smartphones," in *2020 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, pp. 10888–10893, IEEE, 2020.

[51] S. Sun, D. Melamed, and K. Kitani, "Idol: Inertial deep orientation-estimation and localization," in *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 35, pp. 6128–6137, 2021.

[52] O. Asraf, F. Shama, and I. Klein, "Pdrnet: A deep-learning pedestrian dead reckoning framework," *IEEE Sensors Journal*, vol. 22, no. 6, pp. 4932–4939, 2021.

[53] R. Buchanan, M. Camurri, F. Dellaert, and M. Fallon, "Learning inertial odometry for dynamic legged robot state estimation," in *Conference on Robot Learning*, pp. 1575–1584, PMLR, 2022.

[54] M. Zhang, M. Zhang, Y. Chen, and M. Li, "Imu data processing for inertial aided navigation: A recurrent neural network based approach," in *2021 IEEE International Conference on Robotics and Automation (ICRA)*, pp. 3992–3998, IEEE, 2021.

[55] J. Gong, X. Zhang, Y. Huang, J. Ren, and Y. Zhang, "Robust inertial motion tracking through deep sensor fusion across smart earbuds and smartphone," *Proceedings of the ACM on Interactive, Mobile, Wearable and Ubiquitous Technologies*, vol. 5, no. 2, pp. 1–26, 2021.

[56] S. Herath, D. Caruso, C. Liu, Y. Chen, and Y. Furukawa, "Neural inertial localization," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 6604–6613, 2022.

[57] X. Cao, C. Zhou, D. Zeng, and Y. Wang, "Rio: Rotation-equivariance supervised learning of robust inertial odometry," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 6614–6623, 2022.

[58] Y. Wang, J. Kuang, Y. Li, and X. Niu, "Magnetic field-enhanced learning-based inertial odometry for indoor pedestrian," *IEEE Transactions on Instrumentation and Measurement*, vol. 71, pp. 1–13, 2022.

[59] S. S. Saha, S. S. Sandha, L. A. Garcia, and M. Srivastava, "Tinyodom: Hardware-aware efficient neural inertial navigation," *Proceedings of the ACM on Interactive, Mobile, Wearable and Ubiquitous Technologies*, vol. 6, no. 2, pp. 1–32, 2022.

[60] B. Rao, E. Kazemi, Y. Ding, D. M. Shila, F. M. Tucker, and L. Wang, "Ctin: Robust contextual transformer network for inertial navigation," in *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 36, pp. 5413–5421, 2022.

[61] B. Zhou, Z. Gu, F. Gu, P. Wu, C. Yang, X. Liu, L. Li, Y. Li, and Q. Li, "Deepvip: Deep learning-based vehicle indoor positioning using smartphones," *IEEE Transactions on Vehicular Technology*, 2022.

[62] F. Bo, J. Li, and W. Wang, "Mode-independent stride length estimation with imus in smartphones," *IEEE Sensors Journal*, vol. 22, no. 6, pp. 5824–5833, 2022.

[63] H. Tang, X. Niu, T. Zhang, Y. Li, and J. Liu, "Odonet: Untethered speed aiding for vehicle navigation without hardware wheeled odometer," *IEEE Sensors Journal*, 2022.

[64] Y. Wang, H. Cheng, and M. Q.-H. Meng, "A2dio: Attention-driven deep inertial odometry for pedestrian localization based on 6d imu," in *2022 International Conference on Robotics and Automation (ICRA)*, pp. 819–825, IEEE, 2022.

[65] Y. Wang, J. Kuang, X. Niu, and J. Liu, "Llio: Lightweight learned inertial odometer," *IEEE Internet of Things Journal*, 2022.

[66] S. Hochreiter and J. Schmidhuber, "Long short-term memory," *Neural computation*, vol. 9, no. 8, pp. 1735–1780, 1997.

[67] F. Yu and V. Koltun, "Multi-scale context aggregation by dilated convolutions," *The International Conference on Learning Representations (ICLR)*, 2016.

[68] R. S. Sutton and A. G. Barto, *Reinforcement learning: An introduction*. MIT press, 2018.

[69] I. Goodfellow, J. Pouget-Abadie, M. Mirza, B. Xu, D. Warde-Farley, S. Ozair, A. Courville, and Y. Bengio, "Generative adversarial networks," *Communications of the ACM*, vol. 63, no. 11, pp. 139–144, 2020.

[70] E. Tzeng, J. Hoffman, K. Saenko, and T. Darrell, "Adversarial discriminative domain adaptation," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 7167–7176, 2017.

[71] C. Chen, X. Lu, J. Wahlstrom, A. Markham, and N. Trigoni, "Deep neural network based inertial odometry using low-cost inertial measurement units," *IEEE Transactions on Mobile Computing*, 2021.

[72] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, Ł. Kaiser, and I. Polosukhin, "Attention is all you need," *Advances in neural information processing systems*, vol. 30, 2017.

[73] N. A. Abiad, Y. Kone, V. Renaudin, and T. Robert, "Smartstep: A robust step detection method based on smartphone inertial signals driven by gait learning," *IEEE Sensors Journal*, vol. 22, pp. 12288–12297, 6 2022.

[74] M. Jaderberg, K. Simonyan, A. Zisserman, *et al.*, "Spatial transformer networks," *Advances in neural information processing systems*, vol. 28, 2015.

[75] A. Manos, T. Hazan, and I. Klein, "Walking direction estimation using smartphone sensors: A deep network-based framework," *IEEE Transactions on Instrumentation and Measurement*, vol. 71, 2022.

[76] C. Lea, M. D. Flynn, R. Vidal, A. Reiter, and G. D. Hager, "Temporal convolutional networks for action segmentation and detection," in *proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 156–165, 2017.

[77] P. Ren, Y. Xiao, X. Chang, P.-Y. Huang, Z. Li, X. Chen, and X. Wang, "A comprehensive survey of neural architecture search: Challenges and solutions," *ACM Computing Surveys (CSUR)*, vol. 54, no. 4, pp. 1–34, 2021.

[78] R. Clark, S. Wang, H. Wen, A. Markham, and N. Trigoni, "VINet : Visual-Inertial Odometry as a Sequence-to-Sequence Learning Problem," in *The Conference on Artificial Intelligence (AAAI)*, pp. 3995–4001, 2017.

[79] E. J. Shamwell, K. Lindgren, S. Leung, and W. D. Nothwang, "Unsupervised deep visual-inertial odometry with online error correction for rgb-d imagery," *IEEE transactions on pattern analysis and machine intelligence*, 2019.

[80] C. Chen, S. Rosa, Y. Miao, C. X. Lu, W. Wu, A. Markham, and N. Trigoni, "Selective sensor fusion for neural visual-inertial odometry," in *IEEE/CVF International Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 10542–10551, 2019.

[81] L. Han, Y. Lin, G. Du, and S. Lian, "Deepvio: Self-supervised deep learning of monocular visual inertial odometry using 3d geometric constraints," in *2019 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, pp. 6906–6913, IEEE, 2019.

[82] M. R. U. Saputra, P. P. de Gusmao, C. X. Lu, Y. Almalioglu, S. Rosa, C. Chen, J. Wahlström, W. Wang, A. Markham, and N. Trigoni, "Deeptio: A deep thermal-inertial odometry with visual hallucination," *IEEE Robotics and Automation Letters*, vol. 5, no. 2, pp. 1672–1679, 2020.

[83] C. X. Lu, M. R. U. Saputra, P. Zhao, Y. Almalioglu, P. P. De Gusmao, C. Chen, K. Sun, N. Trigoni, and A. Markham, "milliego: single-chip mmwave radar aided egomotion estimation via deep sensor fusion," in

*Proceedings of the 18th Conference on Embedded Networked Sensor Systems*, pp. 109–122, 2020.

[84] P. Wei, G. Hua, W. Huang, F. Meng, and H. Liu, "Unsupervised monocular visual-inertial odometry network," in *Proceedings of the Twenty-Ninth International Conference on International Joint Conferences on Artificial Intelligence*, pp. 2347–2354, 2021.

[85] C. Chen, C. X. Lu, B. Wang, N. Trigoni, and A. Markham, "Dynanet: Neural kalman dynamical model for motion estimation and prediction," *IEEE Transactions on Neural Networks and Learning Systems*, vol. 32, no. 12, pp. 5479–5491, 2021.

[86] Y. Almalioglu, M. Turan, M. R. U. Saputra, P. P. de Gusmão, A. Markham, and N. Trigoni, "Selfvio: Self-supervised deep monocular visual–inertial odometry and depth estimation," *Neural Networks*, vol. 150, pp. 119–136, 2022.

[87] Y. Tu and J. Xie, "Undeeplio: Unsupervised deep lidar-inertial odometry," in *Asian Conference on Pattern Recognition*, pp. 189–202, Springer, 2022.

[88] D. Ramachandram and G. W. Taylor, "Deep multimodal learning: A survey on recent advances and trends," *IEEE signal processing magazine*, vol. 34, no. 6, pp. 96–108, 2017.

[89] E. S. Jones and S. Soatto, "Visual-inertial navigation, mapping and localization: A scalable real-time causal approach," *The International Journal of Robotics Research*, vol. 30, no. 4, pp. 407–430, 2011.

[90] S. Leutenegger, S. Lynen, M. Bosse, R. Siegwart, and P. Furgale, "Keyframe-based visual–inertial odometry using nonlinear optimization," *The International Journal of Robotics Research*, vol. 34, no. 3, pp. 314–334, 2015.

[91] C. Forster, L. Carlone, F. Dellaert, and D. Scaramuzza, "On-manifold preintegration for real-time visual–inertial odometry," *IEEE Transactions on Robotics*, vol. 33, no. 1, pp. 1–21, 2017.

[92] A. Geiger, P. Lenz, C. Stiller, and R. Urtasun, "Vision meets robotics: The kitti dataset," *The International Journal of Robotics Research*, vol. 32, no. 11, pp. 1231–1237, 2013.

[93] M. Burri, J. Nikolic, P. Gohl, T. Schneider, J. Rehder, S. Omari, M. W. Achtelik, and R. Siegwart, "The euroc micro aerial vehicle datasets," *The International Journal of Robotics Research*, vol. 35, no. 10, pp. 1157–1163, 2016.

[94] S. Ahuja, W. Jirattigalachote, and A. Tosborvorn, "Improving Accuracy of Inertial Measurement Units using Support Vector Regression," tech. rep., 2011.

[95] A. Parate, M. C. Chiu, C. Chadowitz, D. Ganesan, and E. Kalogerakis, "RisQ: Recognizing smoking gestures with inertial sensors on a wristband," in *Annual International Conference on Mobile Systems, Applications, and Services (MobiSys)*, pp. 149–161, 2014.

[96] A. Mannini and A. M. Sabatini, "Machine learning methods for classifying human physical activity from on-body accelerometers," *Sensors*, vol. 10, no. 2, pp. 1154–1175, 2010.

[97] A. Valtazanos, D. Arvind, and S. Ramamoorthy, "Using wearable inertial sensors for posture and position tracking in unconstrained environments through learned translation manifolds," in *2013 ACM/IEEE International Conference on Information Processing in Sensor Networks (IPSN)*, pp. 241–252, IEEE, 2013.

[98] M. Yuwono, S. W. Su, Y. Guo, B. D. Moulton, and H. T. Nguyen, "Unsupervised nonparametric method for gait analysis using a waist-worn inertial sensor," *Applied Soft Computing*, vol. 14, pp. 72–80, 2014.

[99] Y. Huang, M. Kaufmann, E. Aksan, M. J. Black, O. Hilliges, and G. Pons-Moll, "Deep inertial poser: Learning to reconstruct human pose from sparse inertial measurements in real time," *ACM Transactions on Graphics (TOG)*, vol. 37, no. 6, pp. 1–15, 2018.

[100] X. Yi, Y. Zhou, and F. Xu, "Transpose: real-time 3d human translation and pose estimation with six inertial sensors," *ACM Transactions on Graphics (TOG)*, vol. 40, no. 4, pp. 1–13, 2021.

[101] X. Yi, Y. Zhou, M. Habermann, S. Shimada, V. Golyanik, C. Theobalt, and F. Xu, "Physical inertial poser (pip): Physics-aware real-time human motion tracking from sparse inertial sensors," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 13167–13178, 2022.

[102] D. Anguita, A. Ghio, L. Oneto, X. Parra Perez, and J. L. Reyes Ortiz, "A public domain dataset for human activity recognition using smartphones," in *Proceedings of the 21th international European symposium on artificial neural networks, computational intelligence and machine learning*, pp. 437–442, 2013.

[103] G. Chevalier, "Lstms for human activity recognition," 2016.

[104] T. Zebin, P. J. Scully, and K. B. Ozanyan, "Human activity recognition with inertial sensors using a deep learning approach," in *2016 IEEE sensors*, pp. 1–3, IEEE, 2016.

[105] D. Ravi, C. Wong, B. Lo, and G.-Z. Yang, "A deep learning approach to on-node sensor data analytics for mobile or wearable devices," *IEEE journal of biomedical and health informatics*, vol. 21, no. 1, pp. 56–64, 2016.

[106] B. M. Eskofier, S. I. Lee, J.-F. Daneault, F. N. Golabchi, G. Ferreira-Carvalho, G. Vergara-Diaz, S. Sapienza, G. Costante, J. Klucken, T. Kautz, *et al.*, "Recent machine learning advancements in sensor-based mobility analysis: Deep learning for parkinson's disease assessment," in *2016 38th Annual International Conference of the IEEE Engineering in Medicine and Biology Society (EMBC)*, pp. 655–658, IEEE, 2016.

[107] J. Windau and L. Itti, "Inertial-based motion capturing and smart training system," in *2019 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, pp. 4027–4034, IEEE, 2019.

[108] K. Weiss, T. M. Khoshgoftaar, and D. Wang, "A survey of transfer learning," *Journal of Big data*, vol. 3, no. 1, pp. 1–40, 2016.

[109] Y. Tian, C. Sun, B. Poole, D. Krishnan, C. Schmid, and P. Isola, "What makes for good views for contrastive learning?," *Advances in Neural Information Processing Systems*, vol. 33, pp. 6827–6839, 2020.

[110] A. Kendall and Y. Gal, "What uncertainties do we need in bayesian deep learning for computer vision?," *Advances in neural information processing systems*, vol. 30, 2017.

[111] J. Gou, B. Yu, S. J. Maybank, and D. Tao, "Knowledge distillation: A survey," *International Journal of Computer Vision*, vol. 129, no. 6, pp. 1789–1819, 2021.

[112] C. Saharia, W. Chan, S. Saxena, L. Li, J. Whang, E. Denton, S. K. S. Ghasemipour, B. K. Ayan, S. S. Mahdavi, R. G. Lopes, *et al.*, "Photorealistic text-to-image diffusion models with deep language understanding," *Neural Information Processing Systems*, 2022.

[113] C. Finn, P. Abbeel, and S. Levine, "Model-agnostic meta-learning for fast adaptation of deep networks," in *International conference on machine learning*, pp. 1126–1135, PMLR, 2017.

[114] B. Mildenhall, P. P. Srinivasan, M. Tancik, J. T. Barron, R. Ramamoorthi, and R. Ng, "Nerf: Representing scenes as neural radiance fields for view synthesis," *Communications of the ACM*, vol. 65, no. 1, pp. 99–106, 2021.

[115] A. Oord, Y. Li, I. Babuschkin, K. Simonyan, O. Vinyals, K. Kavukcuoglu, G. Driessche, E. Lockhart, L. Cobo, F. Stimberg, *et al.*, "Parallel wavenet: Fast high-fidelity speech synthesis," in *International conference on machine learning*, pp. 3918–3926, PMLR, 2018.

**Changhao Chen** is a Lecturer at College of Intelligence Science and Technology, National University of Defense Technology (China). Before that, he obtained his Ph.D. degree at University of Oxford (UK), M.Eng. degree at National University of Defense Technology (China), and B.Eng. degree at Tongji University (China). His research interest lies in robotics, computer vision and cyberphysical systems.



**Xianfei Pan** received the Ph.D. degree in control science and engineering from the National University of Defense Technology, Changsha, China, in 2008. Currently, he is a professor of the College of Intelligence Science and Technology, National University of Defense Technology. His current research interests include Inertial navigation system and indoor navigation system.