# Person Re-Identification and Tracking

# 1. INTRODUCTION

This project presents an advanced person tracking system that combines the strengths of computer vision and machine learning. It's designed to accurately identify, track, and analyze individuals in video streams using state-of-the-art object detection and tracking algorithms. Each person detected in the video is assigned a unique ID, enabling the system to focus on and track specific individuals upon request. This feature is particularly useful for applications where individual tracking across different frames or scenes is crucial.

## 2. SYSTEM ARCHITECTURE

Our system consists of following modules.

- **Configuration Management (config.py)**
  Manages essential system configurations, including model paths, video source, processing options, and image enhancement settings.

- **API Server (app.py)**
  A RESTful API using FastAPI for video frame processing. It listens for video processing requests and initializes the PersonTracker with configurations.

- **Person Tracking Core (tracker.py)**
  The core module for person detection, tracking, and analysis. It utilizes YOLO for object detection, an ReID model for feature extraction of detected persons, and Kalman Filters for movement prediction.

- **Preprocessing Module (preprocessing.py)**
  Enhances video frames for better detection and tracking accuracy. Techniques include image resizing, padding, color normalization, noise reduction, and histogram equalization.

# 3. TEHCNOLOGIES AND MODELS

## 3.1. YOLOv5 Model (yolov5su(small unet) variant)

**Overview:**
YOLO (You Only Look Once) is a popular real-time object detection system. The YOLOv5 architecture by Ultralytics is designed for efficient and accurate object detection. It features multiple versions (YOLOv5s, m, l, x, and n) to balance computational resource use and detection precision. The architecture employs a backbone for feature extraction and a neck for feature fusion, which are crucial for detecting objects across different scales. YOLOv5 utilizes advancements like Cross Stage Partial networks (CSPNet) and Spatial Pyramid Pooling (SPP) to enhance its feature extraction capabilities. This design facilitates real-time object detection with high accuracy, even on devices with limited computational power.

**The YOLOv5 architecture from Ultralytics includes several key components:**

- **Backbone:** The backbone is responsible for initial feature extraction from input images. It processes the images to create a set of feature maps that encode important information about objects in the images.
- **Neck:** The neck performs additional processing on the feature maps from the backbone. It typically involves feature fusion and refinement, enhancing the ability to detect objects of various sizes.
- **Cross Stage Partial Networks (CSPNet):** CSPNet is an architecture that helps in reducing the computational cost while maintaining accuracy. It splits the feature map into two parts and then merges them through cross-stage hierarchy, which is efficient for gradient flow and learning rich features.
- **Spatial Pyramid Pooling (SPP):** SPP is used to make the network more robust to object scale variations. It pools features at different scales and concatenates them together, ensuring that the network can detect objects of various sizes.

**Functionality:**
Divides the input image into a grid. Each grid cell predicts a certain number of bounding boxes and class probabilities for those boxes. Utilizes a single neural network that runs on the entire image, enabling it to predict objects at multiple scales.

## 3.2. Re-identification Model (osnet_ain_x1_0)

**Overview:**
OSNet (Omni-Scale Network) is designed for the task of person re-identification. The osnet_ain_x1_0 model is a specific architecture within the OSNet family.

**Key Features:**
- Incorporates omni-scale feature learning, meaning it can capture and learn features at multiple scales, crucial for accurately differentiating between individuals.
- The architecture is tailored to focus on learning a diverse set of features that can robustly identify a person across different camera views, poses, and lighting conditions.

**Architecture:**
*Omni-Scale Network (OSNet) for person re-identification is quite technical and involves several advanced concepts in neural network design:*

- **Depthwise Separable Convolutions:** This technique used to reduce the computational complexity and size of neural networks. This method divides a standard convolution into two separate layers: depthwise and pointwise convolutions. The depthwise convolution applies a single filter per input channel, and the pointwise convolution then applies a 1x1 convolution to combine the outputs of the depthwise layer. This approach significantly reduces the number of parameters and the computational cost compared to standard convolutions, making the network more efficient while still capable of effective feature extraction.

- **Omni-Scale Residual Block:** Central to OSNet, The Omni-Scale Residual Block in OSNet is a sophisticated component designed to process features at multiple scales. This block is an advanced version of the traditional residual bottleneck, augmented with the Lite 3x3 layer, allowing the network to efficiently capture and integrate features from different scales. It uses stacked Lite 3x3 layers to vary the receptive field, providing a mechanism to handle both small-scale and large-scale features. This design is crucial for the complex task of person re-identification, where identifying individuals across varying conditions and scales is essential.

- **Dynamic and Unified Aggregation Gate:** The Dynamic and Unified Aggregation Gate (AG) in OSNet is a key feature that dynamically fuses multi-scale features from different streams within the omni-scale residual block. It's a neural network-based mechanism that assigns varying weights to different scales based on the input image. This adaptability allows the model to focus on the most relevant scales for a given input, enhancing its efficiency and effectiveness. The AG's design results in a fine-grained fusion of features, tuning each feature channel for optimal performance in person re-identification tasks.

- **Instance Normalization Layers:** Instance Normalization (IN) layers in the OSNet architecture are crucial for enhancing its generalizability across different datasets. Unlike Batch

Normalization, which normalizes based on mini-batch statistics, IN normalizes each sample using its own mean and standard deviation. This approach effectively reduces style differences caused by variations in environments, lighting conditions, and camera setups in different datasets. IN layers thus enable the network to adapt to and perform effectively on diverse datasets, making it particularly suitable for person re-identification tasks that require robustness to varied input conditions.

- **Architecture Search for Optimal IN Placement:** Using neural architecture search (NAS) optimally placing Instance Normalization (IN) layers within the OSNet architecture. This approach involves defining a search space with various candidate blocks incorporating IN in different ways. The objective is to minimize the expectation through joint optimization of the model architecture and parameters. The process involves continuous relaxation and reparameterization strategies to overcome the non-differentiability of discrete selection in architecture search. The outcome of this search is a network architecture with IN layers strategically placed for maximum efficacy in enhancing generalization across different datasets.

# 3.3. Kalman Filters

**Overview:**

The Kalman filter is a powerful tool used in object tracking. It functions by predicting the future state of an object (like its position and velocity) based on its current state and measurement uncertainties. This process involves two main stages: prediction and correction. The filter first makes a prediction about the object's next state, then updates this prediction using new measurements. This cycle allows the filter to refine its tracking accuracy over time, effectively handling uncertainties and noise in the data. filter calculates tracking accuracy by continuously updating its predictions based on incoming data. It compares predicted states (like position and velocity) with actual measurements, and adjusts its predictions to reduce the error. This is achieved through a correction process that involves a mathematical calculation called the Kalman gain. This gain determines how much weight to give new measurements versus the predictions. Over time, this process allows the filter to improve its accuracy in estimating the object's state, even with noisy or incomplete data.

**Example:**

Imagine we're using a drone to track a moving car for a film shoot. The car's exact path is unpredictable due to traffic and turns. The Kalman filter in our drone's tracking system predicts where the car will be based on its last known position and speed. As the drone flies, it constantly receives real-time GPS data on the car's position. The filter compares its predictions with these updates. If there's a difference, it adjusts the prediction to be more accurate, considering both the new data and the previous prediction's uncertainty. This ongoing process allows the drone to accurately follow the car, smoothly adjusting its path and keeping the car in frame despite the uncertain, changing environment.

Real-time updates, like GPS data, provide the current location of the car, but they might not always be perfectly accurate or available at every single moment. The Kalman filter uses these updates along with its predictions to improve tracking accuracy. It's particularly useful when GPS data is noisy, incomplete, or delayed. By continuously adjusting its predictions using new data, the filter helps maintain a more consistent and accurate tracking of the car, ensuring the drone follows it smoothly even when GPS updates are imperfect.

**Our Kalman Filter:**

Captures the essential elements of a Kalman filter for object tracking, as described in the explanation.

- It assumes a constant velocity model and focuses on position estimation.
- "dim" Represents the 4 state variables ($x$, $y$, $\dot{x}$, $\dot{y}$), matching the explanation's focus on position and velocity.
- "F" - Transition Matrix, models how the state evolves over time, assuming constant velocity. It's consistent with the explanation's equations

- "H" - Measurement Matrix, Maps the state to the measured values (x and y coordinates), aligning with the explanation's focus on observing position
- "P" - Initial state uncertainty, initialized with a relatively large value (1000) to reflect potential uncertainty at the start.
- "R" - Measurement noise covariance, representing uncertainty in the observations.
- "Q" - Process noise covariance, capturing uncertainty in the state transitions. It's set using a function that models discrete white noise, matching the explanation's assumptions.
- Tuning the covariance matrices (P, R, Q) might need adjustments based on our application's noise characteristics and desired estimation accuracy.

# PREPROCESSING

The *preprocess_image* function in the *PreprocessImage* class is designed to standardize the input images for further processing in a machine learning pipeline. It begins by resizing the image to maintain the aspect ratio and then adds padding to match a specific target size, ensuring consistency in dimensions across different images. This step is crucial for neural network models that require fixed-size inputs. The function also includes an optional line for histogram equalization, which, if enabled, would enhance the contrast of the image, although it's currently not active. Finally, the image undergoes normalization, converting pixel values from the standard 0-255 range to a 0-1 scale, a common practice for preparing data for deep learning models. This normalization, combined with the resizing and padding, makes the images ready for efficient and consistent processing by subsequent neural network models.

# ANALYSIS
## For RGB Videos:
The person tracking system, in its current form, provides a good starting point for RGB video analysis. While it doesn't yet exhibit excellent performance, it lays a solid foundation for further improvement. The effectiveness of the system, especially in real-world scenarios, can be significantly enhanced by fine-tuning the models with more relevant, real-time filtered data. This process involves adapting the models to be more attuned to the specific characteristics and variations present in real-world environments.

## Fine-Tuning the osnet_ain_x1_0 Model:
### Type of Data Required:
***Diverse Human Images:***
>The data should include a wide range of human subjects, captured in various poses, angles, and lighting conditions. Diversity in clothing, body shapes, and backgrounds is crucial. Real-World Scenarios: Images that reflect real-world situations, like crowded scenes, different urban and indoor environments, are essential to mimic actual operational conditions. Temporal Data: Sequences of images where the same individuals appear in different frames can help the model learn to re-identify people over time and different camera angles.

***Quantity of Data:***
>The amount of data required for effective fine-tuning can vary, but generally, larger datasets lead to better generalization. A dataset comprising thousands of labeled images is a good starting point. For instance, a dataset with 10,000 to 20,000 images could be considered sufficient for initial fine-tuning efforts. The key is not just quantity but also the quality and variety of the data.

### Fine-Tuning Process:
***Data Preparation:***
>This includes labeling the data (if not already labeled), augmenting the dataset (e.g., through random cropping, rotations, scaling), and splitting it into training, validation, and test sets.

***Model Training:***
>Load the pre-trained osnet_ain_x1_0 model and train it on the new dataset. Fine-tuning involves training the model for additional epochs on the new dataset, allowing the model to adjust its weights specifically to the new data characteristics.

***Hyperparameter Tuning:***
>Adjust learning rate, batch size, and other relevant parameters for optimal learning on the new dataset.

*Regular Evaluation:*

Regularly evaluate the model on the validation set to monitor its performance and prevent overfitting.

*Post Fine-Tuning:*

Performance Testing: After fine-tuning, thoroughly test the model's performance on unseen test data to ensure that it has learned to generalize well.

Continuous Improvement: Fine-tuning is an iterative process. Based on the test results, further adjustments in data or training strategy might be necessary.

In conclusion, fine-tuning the osnet_ain_x1_0 model for enhanced performance in the person tracking system involves careful consideration of the dataset's diversity, size, and relevance to real-world conditions. This process, while resource-intensive, is essential for improving the system's accuracy and robustness in various operational scenarios.

# For Thermal Videos:

Improving the person tracking system's performance, particularly for thermal videos, involves addressing the unique challenges posed by thermal imagery. An effective approach to enhance both RGB and thermal video analysis could be to employ a Siamese network architecture based on the osnet_ain_x1_0 model, utilizing NTXent Loss.

**Implementing Siamese Network with osnet_ain_x1_0 and NTXent Loss:**

**Siamese Network Overview:**

A Siamese network consists of twin networks which accept distinct inputs but are joined by an energy function at the top. This setup is powerful for learning embeddings or features where similarity is crucial. In the context of person tracking, this architecture can learn to map the input thermal images into a space where the distance between similar images is minimized, and that between dissimilar images is maximized.

**Utilizing osnet_ain_x1_0:**

The osnet_ain_x1_0 model, known for its effective feature extraction in person re-identification tasks, serves as the backbone of each branch of the Siamese network. This allows the network to leverage osnet_ain_x1_0's ability to discern fine-grained features, crucial in differentiating between individuals in thermal imagery.

**NTXent Loss (Normalized Temperature-scaled Cross Entropy Loss):**

NTXent Loss is particularly suited for learning from pairs of examples (like in Siamese networks), making it ideal for tasks that involve measuring similarity or dissimilarity. In thermal video analysis, this loss function can be instrumental in effectively training the network to understand and encode the nuances of thermal images for accurate person re-identification.

**Data Requirements and Fine-Tuning:**

The dataset for training this Siamese network should include pairs of thermal images of the same and different individuals, under varied environmental conditions and different time frames.

Fine-tuning the network involves training it with this dataset, using NTXent Loss to refine the feature extraction capabilities specifically for thermal imagery. Given the distinct nature of thermal data, the amount of data and diversity in the dataset are crucial for the network to learn effectively.

**Evaluation and Iteration:**

Post fine-tuning, the model's performance should be rigorously evaluated using a separate thermal image dataset. Iterative refinement may be necessary, focusing on improving the model's ability to distinguish between different individuals accurately in thermal images.

In conclusion, integrating a Siamese network structure with the osnet_ain_x1_0 model and employing NTXent Loss presents a promising approach to significantly enhance the person tracking system's capability, especially in handling thermal videos. This method not only leverages the strength of osnet_ain_x1_0 in feature extraction but also capitalizes on the Siamese network's ability to learn robust, discriminative features essential for accurate person re-identification. The success of this approach hinges on careful dataset preparation and continuous model refinement, tailored specifically for thermal video analytics.

## Conclusion

This project demonstrates the integration of deep learning and computer vision for real-world applications, providing a robust platform for person tracking with high configurability and adaptability.