

## Fundamentals

**Computer Vision:** The inverse algorithmic process of inferring world/sensor properties from observations. light source  $\rightarrow$  optical medium  $\rightarrow$  optical sensor + object  $\rightarrow$  observation  $\rightarrow$  CV algorithm

**Dynamic Driving Task (DDT):** Lateral / longitudinal control, object/event detection, recognition, classification, response (OEDR), maneuver planning, enhancing conspicuity via lighting, horn, signaling. NOT trip planning.  
 ↳ Perception, Prediction, Planning, Control

**Driving automation System:** Hardware/software performing parts of DDT.

**Automated Driving System (ADS):** Hardware/software performing whole DDT.

**Operational Design Domain (ODD):** Specific conditions for ADS function.

**Autonomous Car (ADS-DV):** Driverless operation within ODD.

**SAE levels of automation:** 0  $\rightarrow$  none, 1  $\rightarrow$  driver assistance, 2  $\rightarrow$  partial driving automation, 3  $\rightarrow$  "eyes off" but human fallback, 4  $\rightarrow$  high driving automation with limited ODD, 5  $\rightarrow$  full, unlimited ODD

## SENSORS

**External:** Camera, LiDAR, Radar, Event Camera, GPS/IMU

**Internal:** Camera, Thermometer

**Light Detection and Ranging (LiDAR):** active, near-infrared, direct range via time-of-flight.

**Pulsed LiDAR:** Precision  $\approx 7.5\text{cm}$ . Range depends on signal-to-noise ratio, today  $\approx 75\text{cm}$ . Max range  $\propto \frac{1}{\text{emission frequency}}$

$$P_R(R) = G_A \int_0^{2\pi} P_E(\theta) + \left( R - \frac{\sigma}{2} \right) d\theta$$

**Point Cloud acquisition by Scanning spherical coordinates**

$$\begin{pmatrix} x \\ y \\ z \end{pmatrix} = \begin{pmatrix} r \cos \theta \cos \phi \\ r \sin \theta \cos \phi \\ r \sin \phi \end{pmatrix}$$

**Motorized Optomechanical Scanner or MEMS**

$\hookrightarrow$  360° horizontal FOV, but low frame rate and mechanical wear  $\hookrightarrow$  small size, no wear, higher rate but <120° FOV.

**Challenges:** Motion Distortion, non-Lambertian surfaces, optical attenuation, back-scattering, interference

**Camera:**

$$\text{Pinhole model: } \begin{pmatrix} u \\ v \\ 1 \end{pmatrix} = \frac{1}{z_c} K [R | t] \begin{pmatrix} x \\ y \\ 1 \end{pmatrix} \text{ for intrinsics } K = \begin{bmatrix} f_x & s_y & c_x \\ 0 & f_y & c_y \\ 0 & 0 & 1 \end{bmatrix}$$

**Radar:** 70-80 GHz spectrum robust to fog, rain, snowfall.

$$\text{FMCW Scanning radar has Prec} = \frac{P_{\text{trans}} G_{\text{eff}} \lambda^2}{16 \pi^2 r^2 L}$$

Modulate frequency and look at intermediate frequency signal

**Event Camera:** Asynchronous events at single pixels with low latency

$$\Delta L(x_k, t_k) := L(x_k, t_k) - L(x_k, t_k - \Delta t_k) \geq \pm C \text{ irradiance change}$$

## Deep Neural Networks

Activation of neuron  $j$  in layer  $l$ :  $z_{l,j} = \sum_{i=1}^{d_{l-1}} w_{j,i}^{(l)} + \sum_{i=1}^{d_l} w_{j,i}^{(l)} g(z_{l-1,i})$   
 where width of layer  $l$  is  $d_l$  and  $g$  is nonlinear activation func:

$$\text{Sigmoid } g(z) = \frac{1}{1+e^{-z}} \quad \text{Tanh } g(z) = \frac{e^z - e^{-z}}{e^z + e^{-z}} \quad \text{ReLU } g(z) = \max(g(z))$$

$$\check{g}(z) = g(z)(1-g(z)) \quad \tilde{g}(z) = 1-g(z)^2 \quad \check{g}(z) = [z>0]$$

**Binary Classification:** Predict label  $y \in \{0,1\}$  from logit  $NN(x) = z$ .

$$\text{Assuming } \log \hat{P}(y|x) = yz \text{ yields } P(y|x) = \frac{\exp(yz)}{\exp(-z)+\exp(z)} = \text{Sigmoid}(yz).$$

$$\text{Minimize } -\log P(y|x) = \log(1+\exp(-yz)) = \text{Softplus}(-yz).$$

**Multiclass Classification:** Predict label  $y \in [n]$  from logits  $NN(x) = z$ .

$$\text{Assuming } \log \hat{P}(y=i|x) = z_i \text{ yields } P(y=i|x) = \frac{\exp(z_i)}{\sum_j \exp(z_j)} = \text{softmax}(z)_i.$$

$$\text{Minimize } -\log P(y=i|x) = \log(\sum_j \exp(z_j)) - z_i.$$

**Supervised Learning:** Minimize empirical risk  $J(\theta) = \frac{1}{N} \sum_{i=1}^N L(P_\theta(x_i), y_i)$ , usually via SGD. Batch size  $n$  yields standard error  $SD/\sqrt{n}$ . Updates with momentum  $V \leftarrow \alpha V - (1-\alpha)\eta_i \nabla_\theta J_\theta(\theta)$ ,  $\theta \leftarrow \theta + V$  and learning rate  $\eta_t = \eta_0 / (1+t \cdot \text{lr\_rate})$  or  $\eta_t = \eta_0 \exp(-kt)$ .

**Vanishing/exploding gradient solutions:**

Weights are initialized randomly around zero with larger means deeper in the network to avoid vanishing/exploding gradients (clipping). Input is normalized: For channel  $j$ ,  $\hat{x}_{ij} = \frac{x_{ij} - \mu_j}{\sigma_j}$ .

**Pre-activation batch normalization:**

$$\text{Learn } \gamma \text{ and } \beta: \text{For batch } B, \hat{x}_c = \gamma \frac{x_c - \mu_B}{\sqrt{\sigma_B^2 + \epsilon}} + \beta$$

$$\text{Prepare for inference: } \gamma \leftarrow \frac{\mu}{\sqrt{\sigma^2 + \epsilon}}, \beta \leftarrow \beta - \frac{\gamma \mu}{\sqrt{\sigma^2 + \epsilon}}$$

**Regularization:** Early stopping when validation loss is low, weight-lagged penalty, data augmentation (rotation, shearing, warping, cropping, photometric)

## Convolutional Neural Networks

**Convolutional layer:** A kernel has size  $d_{l-1} \times M_l \times N_l$  and produces a 2D feature map by summing the per-channel convolution. Then, add  $d_l$  biases to the result. So,  $d_l(d_{l-1}M_lN_l + 1)$  parameters.

$$z_{l,c}[p,q] = \sum_{k=0}^{d_{l-1}-1} \sum_{i=0}^{M_l-1} \sum_{j=0}^{N_l-1} W_{k,c}[i,j] g(z_{l-1,k}) [sp-1 + r(i-1), sq-1 + r(j-1)] + bc$$

stride  $s_l$  increases receptive field of subsequent neurons non-uniformly, dilation  $r_l$  uniformly sparsely.  $H_l = \left\lfloor \frac{H_{l-1} + 2P_{l-1} - l_1 M_l}{s_l} \right\rfloor + 1$

**Padding layer:** Slide window over image and get min/max/mean.

**AlexNet Architecture:**



## Semantic Segmentation

For inputs  $x \in \mathbb{R}^{H \times W \times 3}$  compute  $\hat{Y} \in \mathbb{R}^{H \times W \times C}$  where  $Y[h,w,c]$  is the softmax-probability of  $x[h,w]$  belonging to class  $c$ . Optimize logits by minimizing  $L_{CE}(\hat{Y}; Y) = -\frac{1}{HW} \sum_{h=0}^H \sum_{w=0}^W \sum_{c=0}^C Y[h,w,c] \log(\hat{Y}[h,w,c])$ .

**One-hot encode  $Y$ :**

$$\text{Mean Intersection over Union: } mIoU(\hat{Y}; Y) = \frac{1}{C} \sum_{c=0}^C \sum_{h=0}^H \sum_{w=0}^W \sum_{c=0}^C \delta[\hat{Y}[h,w]=c] \delta[Y[h,w]=c]$$

Evaluation: Cityscapes, Mapillary Vistas, foggy Cityscapes, ACDC, GTAE

**Bad Approaches:** Sliding window (expensive, no global context), CNN without striding, dilation, or pooling (overfitting)

**Encoder-Decoder via Upscaling:**



**UNet:** Encoder-Decoder with upscaling and concatenation of encoder features in decoder at the same size.

**Transformer:** Allow model to learn receptive field.

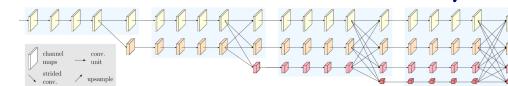
$$\text{Attention } (Q, K, V) = \text{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right)V$$

**Multithead  $(Q, K, V)$ :**  $= \text{Concat}(\text{Attention}(QW_i^Q, KW_i^K, VW_i^V))W^0$

Complexity is  $O(BN^2D)$  in time and  $O(BN^2 + BN^2)$  in space

**Swin Transformer:** Shifted Window Attention as drop-in replacement for convolutional layers. Cross-window connections arise from offsetting windows in deeper layers.

**HRNet:** Parallel multi-resolution convolutions



## Depth Estimation

Monocular from shading, texture, focus (, motion) can be framed as regression task. Use  $L_1, L_2$  loss or Bertluis( $s$ ) =  $\frac{s_{gt}}{\sqrt{s_{pred}^2 + s_{gt}^2}}$  otherwise. Error metrics are rel. difference or RMSE =  $\sqrt{\frac{1}{N} \sum_{i=1}^N (s_{gt} - s_{pred})^2}$ .

**Deep Ordinal Regression:** Discretize depth at points  $t_0, t_1, \dots, t_m$  and predict  $P^k = P[t^k > k]$ . Pixel-wise loss is  $\sum_{i=1}^m (\log(P^k) + \frac{1}{1-P^k})$ . Inferred depth is  $\hat{d} = \frac{t_0 + t_1 + \dots + t_m}{m}$  for  $\hat{t} = \sum_{k=0}^m P^k \geq 0.5$ .

**Plane Coefficient Representation:**  $Z^T = \frac{-a}{fr} u + \frac{-b}{fr} v + \frac{1}{fr} u_0 + \frac{b}{fr} v_0 - c$   
 Normalize by dividing with  $\sqrt{a^2 + b^2 + 1}$ . Then, learn  $(d, \beta, \gamma, \rho)$ .

**P3 Depth:** Learn offset  $\sigma$  st. pto have the same plane coefficients.

**Stereo Depth from disparity:**  $u' = ut_b$ . Then,  $x_w = b \frac{u}{u+b}$ ,  $y_w = b \frac{v}{u+b}$ ,  $z_w = b \frac{u-v}{u+b}$ . Increasing  $b$  and  $\sigma$  increases depth resolution but reduces overlapping FOV. Algorithms match and aggregate cost, optimize weights, then computes 4D cost volume ( $h, w, \text{disparity}, z$ , features).

Regress disparity by linear combination of softmax bucket scores. **Unsupervised stereo depth:** Train monocular estimator and verify by warping and comparing to true stereo image (left-right consistency) OR jointly optimize pose-from-motion and depth prediction from video.

## Object Detection

Detect bounding boxes (or masks for instance segmentation) and classes. Common data sets: PASCAL VOC (small), COCO (large)

**Average Precision (AP):** Precision  $p = \frac{TP}{TP+FP}$ , Recall  $r = \frac{TP}{TP+FN}$

$AP = \int_0^1 AP(r) dr$  where  $pr(r) = \max_{i \in AP} pr_i(r)$  is decreasing stepwise.

$MAP = \frac{1}{C} \sum_{c=0}^{C-1} \sum_{iou=0.5}^1 AP(iou)$ . Pascal VOC uses 11 point avg.

**Slow R-CNN:** Region Proposal  $\rightarrow$  Crop & Warp  $\rightarrow$  CNN  $\rightarrow$  Linear classifier  $\rightarrow$  Task heads

**SPP-net:** Large CNN  $\rightarrow$  ROI Pooling  $\rightarrow$  Small CNN  $\rightarrow$  Task heads

**Multi-Task loss:**  $L(p, u, t^u, v) = -\log p_u + \lambda [u > 1] \sum_{i \in \text{labels}} \text{Huber}(t^u_i - v_i)$ , where  $\text{Huber}(x) = \frac{1}{2}x^2$  for  $|x| \leq 1$  and  $L_{cls}$  not 06  $L_{loc}$

**Faster R-CNN:** Large CNN  $\rightarrow$  Region Proposal  $\rightarrow$  ROI Pooling  $\rightarrow$  MLP  $\rightarrow$  Task heads

**RPN:** A small MLP that is applied Sliding-window style over  $k$  crops to predict objectness & box offsets.

**label objectness truth**  $\hat{p}_i = \begin{cases} 1 & \text{if } i \text{ box. IoU} > 0.7 \\ 0 & \text{if } i \text{ box. IoU} < 0.3 \\ \text{None otherwise} \end{cases}$

$h_{\text{AP}}(P_t) = \frac{1}{N_{\text{cls}}} \sum_i h_{\text{cls}}(A_i, p_i^*) + \frac{1}{N_{\text{reg}}} \sum_i p_i^* h_{\text{reg}}(t_i, p_i)$

**Feature Pyramid Network (FPN):** Append prediction head to U-Net decoder layers.

**Direct set prediction for Object Detection:** Instead of non-maximum suppression, compute bipartite matching via Hungarian algorithm.

**Detection Transformer (DETR):** Encoder-decoder attention and  $N$  object queries which are processed one-by-one by a Feed-Forward Network.

## Instance Segmentation

**Mask R-CNN:** Detect-then-segment, so ...  $\rightarrow$  ROIAlign  $\rightarrow$  MLP  $\rightarrow$  Box regressor  $\rightarrow$  classifier  
**ROIAlign** performs interpolation. Masks are upsampled and unthresholded.  
 $\Theta$  low-res, expensive inference, near-duplicate predictions

**Spatial Mappings:** Proposal-free method.  
 $\text{CNN} \rightarrow \text{seed CNN} \sim \text{per-class seed map} \rightarrow \text{Sample Instance CNN} \sim \nabla\text{-maps for size and pixel offsets for instance} \rightarrow \text{Clusters} \phi_k(\mathbf{c}_i) = \exp(-\| \mathbf{l} - \mathbf{c}_i \|_2^2 / (2\sigma_k^2))$

**Panoptic Segmentation:** Countable instances AND uncountable classes without overlap. Predict  $Y_i \in \mathcal{E}_1, \dots, (\mathcal{C} \times \mathcal{E}_1, \dots, \mathcal{E}_k)$ .

**Panoptic Quality:** Each GT segment  $p_i = \frac{\text{IoU}(A_i, p_i)}{1 - \text{IoU}(A_i, p_i)}$  can have  $\leq 1$  predicted segment with  $\text{IoU} > 0.5$ . The unique matching yields segmentation quality FT-recognition quality.

**MaskFormer:** Variants of DETR with parallel pixel-level module.  
 $\Theta$  Slow, not sort in instance OR semantic, trained per-task.

**OneFormer:** Provide textual queries for the model to search for:  
 $\text{Segmentation} \rightarrow \text{["car", "person", "road"]}, \text{Panoptic: } 8 \cdot 2 \cdot \text{["car"]} + \text{["road"]}.$

## 3D Sensing

**Multi-View CNN:** Render to 2D and run CNN per-image. Then, pool shape descriptors with another CNN.  $\Theta$  linear in number of views, unclear how noisy/incomplete renders are handled.

**Volumetric CNN:** Work on voxel representation.  $\Theta$  information lost during voxelization, sparsity increases with grid resolution.  
Alternative: Oct Trees for recursive partitioning.

**PointNet:** Process scattered, unordered, noisy point clouds.

**Input Transform**  $\rightarrow$  MLP  $\rightarrow$  Feature Transform  $\rightarrow$  MLP  $\rightarrow$  Max Pooling  $\rightarrow$  Classifier (pointwise)  
Transforms are learned matrices with the goal of finding a canonical representation for geometric invariance. Permutation invariance from Maxpool.  
 $\Theta$  No local context for each point, global features depend on coordinates.

**PointNet++:** Iterative furthest Point Sampling  $\rightarrow$  Grouping via ball query  $\rightarrow$  PointNet  $\rightarrow \dots \rightarrow$  Segmentation Head  
 $N \times d(c) \xrightarrow{\text{N} \times \text{nd centroids}} N \times k \times d(c) \xrightarrow{\text{N} \times d(C)}$   
For segmentation, upscale with skip connection in interpole layer.  $k$ -nearest neighbors' features are weighted by distance:  $w_i = \frac{1}{\|x_i - x_j\|}$

## 3D object Detection

**Birds Eye View (BEV):** 2D top-down IoU.

**VoxelNet:** Object proposals from Voxels.

Partition  $\rightarrow$  Group & Sample  $\rightarrow$  Stacked Feature Encoding on centroid-augmented points  $\rightarrow$  3D-CNN (with Batch Norm & ReLU)  $\rightarrow$  PAF  
Trained via Faster-RCNN loss. For regression, use  $\Delta x = \frac{x_b - x_a}{\sigma^2}, \Delta l = \log(\frac{l_b}{l_a}), \Delta \theta = \theta_b - \theta_a$

**3D Point Cloud Processing Pipeline:** 3D Input  $\rightarrow$  Point Cloud  $\rightarrow$  Intermediate Layer  $\rightarrow$  Point-wise Feature  $\rightarrow$  Point-wise MLP  $\rightarrow$  Point-wise Concentrate

## Point RCNN: 2-stage detector with 3D proposal, then refinement.

**Bin-based 3D box generation:** Establish a grid around each foreground point. Classify bin-x, bin-y, bin-theta and regress z, h, w, l, Ax, Ay, Az.  $\text{bin}_x^{(i)} = \left\lfloor \frac{x^i - x^{(i)}}{\delta} \right\rfloor + S$  with  $S = \frac{25}{H \cdot \text{bins}}$ ,  $\text{res}_x^{(i)} = \frac{1}{\delta} (x^i - x^{(i)} + 5 - (\text{bin}_x^{(i)} \cdot \delta + \frac{\delta}{2}))$ . The per-point loss for bins is  $L_{cls} + L_{reg}$

**3D Localization & Reconstruction**

**Iterative Closest Point (ICP):** Obtain  $\{\hat{R}|t\}$  and assignments  $\hat{\phi}: \text{In} \rightarrow \text{In}'$  such that  $E(R, t, \phi) = \sum_i \|R\hat{p}_i(p_i) + t - p_i^{(t)}\|_2^2$  is minimized. Solved iteratively by guessing, computing  $\phi$  by closeness, then updating  $\{\hat{R}|t\}$  until convergence. closed-form  $\hat{t} = M_3 - R M_2 \phi$  (align contracts), optimal rotation by aligning principal 3D directions of variation. (Split  $R = UV^T$  and maximize brace  $(UWV)$  where  $W = \hat{E}(p_1^{(t)}, \dots, p_N^{(t)}) (p_1^{(t)} \dots p_N^{(t)})^T = \hat{U}EV^T$ ). Eliminating reflections yields  $R = \hat{U}[\hat{p}_1^{(t)}, \dots, \hat{p}_N^{(t)}]V^T$   $\Theta$  sensitive to outliers, expensive

**Visual Odometry:** Estimate relative 3D Euclidean motion of a camera from its captured images. Motion from  $I_{k-1}$  to  $I_k$  is  $[R_{k,k-1} \mid t_{k,k-1}]$ .

**2-frame Structure from Motion:** Epipoles  $e_1$  and  $e_2$  are the intra-plane intersections of the baseline offset  $\overrightarrow{O_1 O_2}$ . Potential matches  $p_2$  to  $p_1$  must lie on the epipolar line  $\overleftrightarrow{e_1 p_2}$ :  $p_2^{(1)}(t \times p_1^{(1)}) = (p_2^{(2)})^T [E]_x R p_1^{(1)} = 0$

Computing the essential matrix  $[E]_x R$  needs 5 point correspondences at minimum. The 8-point algorithm yields 4 solutions which are filtered by triangulation.

**RANSAC:** Randomly select minimal correspondences. Estimate structure from motion. For all other points, compute distance to induced epipolar line. Store inliers. Repeat  $N$  times. Finally, compute with the largest set of inliers.

**Super Glue Feature Matching:** Compute pairwise feature similarity matrix via Attentional Aggregation. Then, run Sinkhorn Algorithm to refine into soft assignment matrix where  $\sum_j p_{ij} = 1$ .

**Neural Radiance Fields (NeRFs):** Learn volume rendering densities from posed images via MLPs. Observed color for ray  $\mathbf{r} = \mathbf{o} + \mathbf{t}\mathbf{d}$  is  $C(\mathbf{r}) = \int T(t) \mathbf{f}(r(t)) C(r(t), d) dt$ . Predict via stratified quadrature:  $\hat{C}(\mathbf{r}) = \sum_i T_i [1 - \exp(-\tau_i \cdot (t_i - t))] c_i$ . Here,  $T(t) = \exp(-\int_t^\infty \mathbf{f}(s) ds)$  is the accumulated transmittance and  $c$  color. To improve resolution, use positional embeddings in  $> 5D$  space for input.

**NeRF in the wild:** Robust learning of static classifies and colors from images taken on different days by Learning per-image transient embeddings

**Block NeRF:** Partition based on Visibility.

## Domain Adaptation

**Domain shift:** Formally, images  $I$  and ground truths  $Y$  are sampled from domains  $s$  and  $t$ :  $(I_s, Y_s) \sim p_s(I, Y)$  and  $(I_t, Y_t) \sim p_t(I, Y)$ . When  $p_s \neq p_t$ , we have domain shift. E.g., weather conditions, sensor setup, real vs. synthetic, object semantics changing in different countries. (low  $\rightarrow$  high)

**Weakly Supervised DA:** Label correspondences between source and target inputs, discarding similar outputs.

**Source-Free DA (pre-test), Online DA (test-time), Domain Generalization (pre-training)**

**Input-level DA:** Foggy Cityscapes (synthetic fog based on depth estimation), LiDAR Snowfall, image translation with CycleGANs, Fourier Domain Adaptation.

## CyCADA: Adversarial loss at feature-level to make differentiation from target domain impossible.

$h_{\text{src,feat}} = [\mathbb{E}_{I \sim s, x \sim I} [\log D_{\text{feat}}(p_f(I_t))]] + [\mathbb{E}_{I \sim s, x \sim I} [\log (1 - D_{\text{feat}}(p_f(G_{\text{src}}(I_s))))]]$

**CISS: Contrastive learning of Siamese networks with shared features.**

**CMA:** Patch-wise learning with contrastive embeddings.

**Output-level DA:**

**AdaptSegNet:** Optimize for quality and indistinguishability in an adversarial setting.  $\Theta$  Biases towards easily transferable classes.

**DAIC:** Mix images from two domains together based on their masks, hoping to eliminate implicit training of per-domain submodels. (e.g. put Mexican Taxis on Berlin streets)

**Probabilistic Warping:** Learn deterministic warps to find correspondances.

**Test-time adaptation via TGAT:** Optimize batch norm statistics and parameters to better reflect the test set by minimizing entropy.

## Multi-Modal Detection

**Sensor Fusion:** Late (fuse results), early (fuse inputs), deep (fuse features), asymmetric, weak.

**Frustum PointNet:** Use camera to constrain 3D search space.

**MV3D:** Generate proposals from BEV for front-view LiDAR, BEV, and camera. Then, deeply fuse sensors.

**Point Painting:** Color point clouds from camera image pre-input.

**TransFusion:** Use attention to let model decide where to look up camera.

**COC:** Late fusion of camera-lidar object candidates.

**CRF-Net:** Camera-radar fusion by fusing height with pillars.

Entropy-based fusion VS. MWC-Attention: Both quadratic in time, former might fail with snow.

## Misc. Improvements

**Visual Grounding:** Given a scene and a verbal query, detect the referred object in the scene. 3DV-T does verbo-visual fusion with transformers.  $\Theta$  Spatial specifications ambiguous ("left of x"). Solved via multi-view encoding (MVT-30VG) or learned camera embeddings (Concrete Net).

**Anomaly Segmentation:** Find areas of episodic (model-related) uncertainty, not aleatory (data-related) uncertainty, by doing Monte-Carlo integration using dropout (find large  $\log(\frac{1}{4} \sum_i Y_i^2) - \frac{1}{4} \sum_i Y_i \log(Y_i)$ ) OR image resynthesis with discrepancy network.

**Vehicle-to-Vehicle Communication for Perception (V2VNet):** Maximize informativeness facing unknown time delays and limited bandwidth

## Object Tracking & Path Prediction

**Tracking by detection:** In both frames, detect objects. For each box, compute features. Associate boxes by feature similarity.

**MOT metrics:** Mostly tracked trajectories (>80% of frames), fragments (>80% of trajectories), 1D switches

$\text{MOTA} = 1 - \frac{\text{FN} + \text{FP} + \text{IDSW}}{\text{GT}} \in [-\infty, 1]$

**MOTR:** End-to-end, MOTR inspired prediction of object sequences.

**CNN-based path prediction:** Encode past frame features as images and perform sensor fusion with them.

**RNN-based path prediction:** Use Long Short-Term Memory encoders.