

Gemini Chat



ยินดีครับ ผมจะสอนเนื้อหาเรื่อง **Object Detection** จากเอกสารที่คุณ
แนบมาให้อ่านย่างละเอียด โดยจะเจาะลึกไปถึงลักษณะเพื่อให้เห็นวิวัฒนาการ
และเหตุผลเบื้องหลังของแต่ละเทคนิคครับ

เราจะเริ่มกันที่ **Part 1: พื้นฐานและปัญหาของ Object Detection** ครับ

Part 1: พื้นฐานและปัญหาของ Object Detection (Fundamentals & Challenges)

ก่อนที่เราจะไปดูโมเดล AI ล้ำๆ เราต้องเข้าใจก่อนว่า "โจทย์" ของเรามีอะไร และทำให้มันถึงยากกว่างานทั่วไปครับ

1. ความแตกต่างระหว่าง Classification และ Detection

ในงาน Computer Vision เรา มีงานหลายประเภทที่คล้ายกันแต่ไม่เหมือนกันครับ:

- Image Classification:** รับภาพเข้ามา ตอบว่า "ภาพนี้คืออะไร" (เช่น แมว) โดยไม่สนใจว่าแมวอยู่ตรงไหน
- Object Detection:** แยกกว่ามาก เพราะต้องตอบ 2 อย่างคือ "**WHAT**" (วัตถุคืออะไร) และ "**WHERE**" (วัตถุอยู่ตรงไหน - โดยระบุเป็นกรอบสี่เหลี่ยม หรือ Bounding Box)

2. ความท้าทายหลัก (The Core Challenges)

ทำไมเราถึงใช้ Neural Network แบบเดิมๆ (Standard CNN) มาทำ Detection เลยไม่ได้?

- Multiple Objects:** ภาพหนึ่งภาพอาจมีวัตถุหลายชิ้น หรือไม่มีเลย เราไม่รู้ล่วงหน้าว่าต้อง Output กี่คำตอบ
- Variable Types:** เราต้องทำทั้งจำแนกประเภท (Classification) และหาพิกัด (Regression) พร้อมกัน

3. **High Resolution:** งาน Detection ต้องใช้ภาพความละเอียดสูงกว่างาน Classification ทั่วไป เพื่อให้เห็นวัตถุเล็กๆ ได้ชัดเจน (เช่น 800x600 เทียบกับ 224x224)

3. แนวคิดดังเดิมที่ "ล้มเหลว": Sliding Window

วิธีที่ตรงไปตรงมาที่สุดคือ **Sliding Window** หรือการเลื่อนกรอบสี่เหลี่ยมไปก้าวภาพครับ

- **วิธีการ:** ตัดภาพส่วนเล็กๆ (Crop) มาเข้า CNN เพื่อถามว่า "ใช่ไหม ไหม?", "ใช่แมวไหม?" ถ้าไม่ใช่ก็เลื่อนกรอบไปเรื่อยๆ
- **ปัญหา:**
 - เราไม่รู้ขนาดของวัตถุ จึงต้องลองกรอบหลายขนาด (Scales) และหลายสัดส่วน (Aspect Ratios)
 - ในภาพขนาด 800x600 หากเราลองทุกความเป็นไปได้ จะมีกรอบให้ตรวจสอบถึง **58 ล้านกรอบ!** (58M boxes)
 - **สรุป:** วิธีนี้ซ้ำเกินไปจนทำจริงไม่ได้ครับ (Computational Expensive)

4. ทางออกแรก: Region Proposals (Selective Search)

เพื่อแก้ปัญหาการภาดทุกจุด นักวิจัยจึงเสนอให้ใช้ **Region Proposals**

- **หลักการ:** แทนที่จะสุมทุกจุด ให้ใช้อัลกอริทึม (ที่ไม่ใช่ AI แต่เป็น Heuristics) เพื่อหาริเวณที่ "น่าจะเป็นวัตถุ" (Blob-like regions)
- **Selective Search:** เป็นอัลกอริทึมยอดนิยมในยุคนี้ โดยจะดู สี, พื้นผิว (Texture), และความต่อเนื่องของภาพ เพื่อเสนอกรอบสี่เหลี่ยมมาประมาณ **2,000 กรอบ** ต่อภาพ
- **ข้อดี:** ลดจำนวนพื้นที่ที่ต้องตรวจสอบจาก "58 ล้าน" เหลือแค่ "2 พัน" ทำให้การใช้ CNN เป็นไปได้มากขึ้น

สรุป Part 1: งาน Object Detection ยาก เพราะต้องหาตำแหน่งวัตถุที่มีจำนวนไม่แน่นอน การใช้ Sliding Window นั้นเป็นไปไม่ได้ในทางปฏิบัติ จึงเกิดแนวคิด **Region Proposals** (เช่น Selective Search) ขึ้นมา

เพื่อคัดกรองพื้นที่ที่น่าสนใจก่อนส่งให้ AI ประมวลผล ซึ่งนี่คือภูมิใจ
สำคัญที่นำไปสู่โมเดล **R-CNN** ครับ

หากคุณพร้อมสำหรับ **Part 2: R-CNN และวิธีการวัดผล** พิมพ์ **next**
ได้เลยครับ

- ◆ เยี่ยมครับ มาต่อ กันที่ **Part 2: R-CNN และการวัดผล (Evaluation)**
ครับ

ในส่วนนี้เราจะเจาะลึกโมเดลแรกที่นำ Deep Learning มาใช้กับ Object
Detection ได้สำเร็จ นั่นคือ **R-CNN (Region-based CNN)** และ
เรียนรู้วิธีวัดความแม่นยำของโมเดลพวณ์ครับ

Part 2: R-CNN (Region-based CNN)

R-CNN คือการเอาแนวคิด Region Proposals จาก Part 1 มาผสานกับ
พลังของ CNN ครับ

1. ขั้นตอนการทำงานของ R-CNN (The Pipeline)

กระบวนการทำงานแบ่งเป็น 4 ขั้นตอนหลักที่กำหนดเป็นลูกโซ่:

- Input & Region Proposals:** รับภาพเข้ามา และใช้ Selective Search หาบริเวณที่น่าสนใจ (Region of Interest - **RoI**)
ประมาณ 2,000 จุด
- Warping (การบิดภาพ):**
 - ปัญหา:** CNN ทั่วไป (เช่น AlexNet, VGG) ต้องการภาพ input
ขนาดต่างๆ (เช่น 224x224) และ RoI ที่ตัดมาจากขนาดไม่เท่ากัน
 - วิธีแก้:** บังคับยืด/หดภาพ (Warp) จาก RoI ให้กลมเป็นขนาด
224x224 เท่ากันหมด เพื่อส่งเข้า CNN
- Feature Extraction:** ส่งภาพ 2,000 ภาพที่ Warp และ เข้าไปใน
CNN เพื่อดึง Feature ออกมา (ได้เฉพาะขนาด 4096)
- Classification & Regression:** นำ Feature ที่ได้ไปทำ 2 อย่าง
พร้อมกัน:
 - Classify:** กำหนดว่า "กล่องนี้คืออะไร" (ใช้ SVM ในยุคแรก)

- **Bbox Regression:** กำหนด "ค่าปรับแก้พิกัด" (Offset) เพื่อขยับกรอบให้แม่นยำขึ้น

2. การปรับแก้กรอบ (Bounding Box Regression)

ทำไมต้องปรับกรอบ? เพราะ Region Proposal ที่ได้มาอาจจะแค่ "เฉียดๆ" วัตถุ แต่ไม่เป็น

- โมเดลจะไม่กายพิกัด (x, y, w, h) ตรงๆ แต่จะกาย **ค่าการแปลง (Transform values)** 4 ตัว คือ (t_x, t_y, t_w, t_h)
- **Translation:** t_x, t_y ใช้เลื่อนจุดศูนย์กลางกล่อง
- **Scale:** t_w, t_h ใช้ยืดหดความกว้างและสูง (ผ่านพังก์ชัน exponent เพื่อกันค่าติดลบ)

3. ตัวชี้วัดความแม่นยำ (Metrics: IoU & mAP)

ในงาน Detection เราวัดแค่ "ถูก/ผิด" ไม่ได้ เพราะกรอบที่กายอาจจะเหลือมกับของจริง เราจึงใช้:

- **IoU (Intersection over Union):** วัดว่ากรอบที่กาย (Prediction) ซ้อนกับกับจริง (Ground Truth) มากแค่ไหน
 - **สูตร:** $\frac{\text{พื้นที่ที่ซ้อนทับกัน}}{\text{พื้นที่รวมของทั้งสองกล่อง}}$
 - **เกณฑ์:** โดยทั่วไป $\text{IoU} > 0.5$ ถือว่า "เจอวัตถุ" (Decent), $\text{IoU} > 0.7$ ถือว่า "ดีมาก", $\text{IoU} > 0.9$ คือ "เกือบสมบูรณ์แบบ"
- **mAP (Mean Average Precision):** เป็นเกรดเฉลี่ยของโมเดล
 1. คำนวนกราฟ Precision-Recall ของแต่ละคลาส (เช่น กราฟของเมว, กราฟของหมา) โดยดูจากค่าความเชื่อมั่น (Confidence Score)
 2. หาพื้นที่ใต้กราฟ (Area Under Curve) ของแต่ละคลาส จะได้ค่า **AP (Average Precision)**
 3. นำ AP ของทุกคลาスマร่วมกันแล้วหาค่าเฉลี่ย จะได้ **mAP**

จุดตายของ R-CNN (The Bottleneck)

ถึงแม้ R-CNN จะแม่นยำ แต่มีปัญหาใหญ่หลังที่ทำให้ "ช้ามาก":

- เรายังต้องรัน CNN ถึง 2,000 ครั้งต่อ 1 ภาพ! (เพราะมี 2,000 Regions ต่างคนต่างรัน)
- การ Process ภาพเดียวอาจใช้เวลาเกือบนาที ซึ่งใช้จริงแบบ Real-time ไม่ได้เลย

นักวิจัยจึงคิดค้นวิธีแก้ปัญหาโดยการ "รัน CNN และครองเดียวพอ" ซึ่งเป็นกีมาของ **Fast R-CNN**

ถ้าพร้อมที่จะไปดูวิธีแก้ปัญหานี้ใน Part 3: Fast R-CNN และ RoI Pooling พิมพ์ **next** ได้เลยครับ

- ◆ มาต่อ กันที่ **Part 3: Fast R-CNN และ RoI Pooling** (การแก้ปัญหาความซ้ำ) ครับ

ในส่วนนี้เราจะมาดูว่า นักวิจัยแก้ปัญหา "ครอบครัว" ของ R-CNN กี่ต้องรัน CNN ช้าๆ 2,000 ครั้งได้อย่างไร จนเกิดเป็นโมเดลที่เร็วกว่าเดิมแบบก้าวกระโดดครับ

Part 3: Fast R-CNN

หัวใจสำคัญของ **Fast R-CNN** คือแนวคิดที่ว่า "อย่าทำซ้ำในสิ่งที่กำไรแล้ว" ครับ

1. การแชร์ฟีเจอร์ (Shared Features)

ใน R-CNN เดิม เราตัดภาพ 2,000 ชิ้นแล้วส่งเข้า CNN ทีละชิ้น ซึ่งส่วนใหญ่เป็นภาพที่ซ้อนกับกัน (Overlapping) ทำให้เราคำนวณซ้ำซ้อนมากคลาด

- **วิธีใหม่:** เราจะส่ง **ภาพทั้งภาพ (Whole Image)** เข้า CNN และ ครองเดียว ໄປเลย
- **ผลลัพธ์:** เราจะได้ **Feature Map** ขนาดใหญ่ที่เป็นตัวแทนของภาพทั้งภาพ ออกมาหนึ่งอัน

2. การฉายภาพกลับ (Projection)

เมื่อเรามี Feature Map ของทั้งภาพแล้ว แต่เรายังต้องการฟีเจอร์เฉพาะจุด (RoI) จาก Selective Search อยู่ เราจะทำอย่างไร?

- เราใช้วิธี **Project** (ฉายภาพ) กรอบสี่เหลี่ยมจากภาพต้นฉบับ ลงไปบน Feature Map โดยตรงครับ
- เนื่องจาก CNN ย่อขนาดภาพลง (Subsampling) เราจึงสามารถคำนวณพิกัดที่สัมพันธ์กันได้ (เช่น ถ้าภาพถูกย่อลง 32 เท่า พิกัด (x, y) บนภาพจริง ก็จะอยู่ที่ $(x/32, y/32)$ บน Feature Map)

3. นวัตกรรมสำคัญ: RoI Pooling

นี่คือประโยชน์ของ Fast R-CNN ครับ เพราะหลังจาก Project กรอบลงมาบน Feature Map และ กรอบแต่ละอันจะมีขนาด ไม่เท่ากัน แต่เลเยอร์ตั้งไป (Fully Connected Layer) ต้องการ input ขนาด คงที่ (เช่น 2×2 หรือ 7×7)

RoI Pooling คือเทคนิคการแปลง Feature ขนาดเท่าไหร่ก็ได้ ให้กลายเป็นขนาดคงที่:

- Snap to Grid:** ปัดเศษพิกัดของกรอบให้ลงล็อกกับตาราง (Grid) ของ Feature Map
- Divide:** แบ่งพื้นที่ในกรอบออกเป็นตารางย่อยๆ ตามขนาดที่ต้องการ (เช่น ต้องการ output 2×2 ก็แบ่งเป็น 4 ช่อง)
- Max Pool:** หาค่าสูงสุดในแต่ละช่องย่อย เพื่อให้ได้ค่าตัวแทน
 - ผลลัพธ์:** ไม่ว่ากรอบต้นฉบับจะกว้างยาวแค่ไหน สุดท้ายจะได้เวกเตอร์ขนาดเท่ากัน เป็น พร้อมส่งไปจำแนกประเภทต่อได้เลย

4. ข้อเสียของ RoI Pooling และทางแก้ (RoI Align)

แม้ RoI Pooling จะดี แต่การ "ปัดเศษ" (Snap to Grid) ทำให้ข้อมูลคลาดเคลื่อนเล็กน้อย (Misalignment) ซึ่งส่งผลกระทบมากถ้าวัดความแม่นยำระดับพิกเซล

- RoI Align:** เป็นเวอร์ชันอัปเกรด (มักใช้ใน Mask R-CNN) โดยเลิกใช้การปัดเศษ แต่ใช้คณิตศาสตร์เรื่อง **Bilinear Interpolation** (การประมาณค่าในช่วง) เพื่อคำนวณค่าจุดทุกชนิดแทน ทำให้ได้ฟีเจอร์ที่แม่นยำกว่ามาก

สรุปความเร็ว (Speedup)

Fast R-CNN เร็วกว่า R-CNN อย่างมหาศาล:

- เวลาเท่านั้น: จาก 84 ชั่วโมง เหลือ 8.75 ชั่วโมง
- เวลาใช้งานจริง (Test time): เร็วกว่าเดิมมาก แต่... ยังติดปัญหาเดิม

ปัญหาใหม่ (The New Bottleneck): แม้ตัว AI จะรันเร็วมาก แต่บันทึกผลการหา **Region Proposals (Selective Search)** ยังเป็นอัลกอริทึมเด่าที่รันบน CPU ซึ่งกินเวลาไปเกือบกึ่งกระบวนการ (ประมาณ 2 วินาทีต่อภาพ ในขณะที่ AI ใช้เวลาเสี้ยววินาที)

เพื่อให้เร็วระดับ Real-time เราต้องกำจัด Selective Search ก็งไปครับ

ถ้าพร้อมที่จะดูว่าเราทำอย่างไรใน **Part 4: Faster R-CNN และ Region Proposal Network (RPN)** พิมพ์ **next** ได้เลยครับ

- ◆ เยี่ยมเลยครับ เรามาต่อ กันที่ **Part 4: Faster R-CNN และ Region Proposal Network (RPN)** ซึ่งเป็นจุดเปลี่ยนสำคัญที่ทำให้ Object Detection เร็วพอที่จะใช้งานจริงได้ครับ

Part 4: Faster R-CNN

ใน Part 3 เราเห็นแล้วว่า Fast R-CNN แก้ปัญหาการรันโมเดลซ้ำซ้อนได้แต่ยังติดปัญหาที่ **Selective Search** (การหา Region Proposals) นั้นช้าและรันบน CPU ไม่ได้เรียนรู้ไปพร้อมกับโมเดล (Not Learnable)

Faster R-CNN จึงเกิดมาเพื่อตอบโจทย์นี้โดยเฉพาะ ด้วยแนวคิด "ใช้ CNN หา Region Proposals เองนะเลย"

1. Region Proposal Network (RPN)

แทนที่จะใช้อัลกอริทึมภายนอก (Selective Search) เราจะสร้าง Neural Network เล็กๆ ขึ้นมาตัวหนึ่งชื่อว่า **RPN** และแทรกมันเข้าไปในโมเดลครับ

- **ตำแหน่ง:** RPN จะอยู่ต่อจาก "Backbone CNN" (หลังจากได้ Feature Map แล้ว) และอยู่ก่อน "Classifier"
- **หน้ากี:** รับ Feature Map เข้ามาแล้วพ่น "กรอบสี่เหลี่ยมที่น่าจะเป็นวัตถุ" (Proposals) ออกมาให้เรา

2. หัวใจสำคัญ: Anchor Boxes

RPN ทำงานอย่างไร? มันจะเลือกหน้าต่างเล็กๆ (Sliding Window) ไปทั่ว Feature Map แต่ละจุด

- ที่แต่ละจุด มันจะจินตนาการ "สมอ" (Anchors) หรือกรอบสี่เหลี่ยมมาตรฐานขึ้นมาหลายๆ แบบพร้อมกัน (เช่น 3 ขนาด x 3 ตัวเลือก = 9 Anchors ต่อจุด)
- ทำไมต้องมีหลายแบบ?: เพื่อให้รองรับวัตถุที่หลากหลาย เช่น คน (ผู้ชาย/ผู้หญิง), รถ (กวางเตี้ย), หรือสุนัข (จักรยาน)

3. การทำงานของ RPN

สำหรับ Anchor แต่ละอันกี่จุดหนึ่งๆ RPN จะคำนวณ 2 อย่าง:

1. **Objectness Score:** คำแนะนำว่า "ในกรอบนี้มีวัตถุไหม?" (ใช่/ไม่ใช่)
2. **Box Regression:** ค่าปรับแก้พิกัด (t_x, t_y, t_w, t_h) เพื่อขยับ Anchor ให้ตรงกับวัตถุจริงมากขึ้น

4. โมเดลแบบ 2 ขั้นตอน (Two-Stage Detector)

Faster R-CNN ถือเป็นต้นแบบของ **Two-Stage Detector** ที่สมบูรณ์แบบ:

- **Stage 1 (RPN):** คาดคะเนว่าภาพมีวัตถุอยู่ที่ใดเพื่อคัดกรอง "ว่าที่วัตถุ" (Proposals) ออกมาจำนวนหนึ่ง (เช่น 300-2,000 กล่อง)
- **Stage 2 (Refinement):** นำกล่องเหล่านั้นไปเข้ากระบวนการ ROI Pooling (เหมือน Fast R-CNN) เพื่อจำแนกประเภทและอ้างอิงตำแหน่งของวัตถุ

5. ผลลัพธ์

- **ความเร็ว:** เร็วขึ้นจนเกือบเป็น Real-time (ประมาณ 5-7 fps ในยุคนั้น) เพราะทุกส่วนรันบน GPU ทั้งหมด
 - **ความแม่นยำ:** สูงขึ้น เพราะตัวหา Region Proposals (RPN) ถูกтренมาเพื่อให้เก่งขึ้นเรื่อยๆ (Learnable) ไม่ได้ใช้กฎตายตัวเหมือน Selective Search
-

คำถามที่นำไปสู่อนาคต

Faster R-CNN ทำงานได้ยอดเยี่ยม แต่การมี 2 ขั้นตอน (Stage 1 หาบอง -> Stage 2 จำแนกของ) ทำให้มันยังมีความซับซ้อนและซ้าอยู่บ้าง

นักวิจัยจึงตั้งคำถามว่า "เราจำเป็นต้องมี 2 ขั้นตอนจริงๆ หรือ? ทำไมไม่ทำนายทุกอย่างออกมากีเดียวเลยล่ะ?"

นี่คือจุดเริ่มต้นของ **Single-Stage Detectors** ที่ดึงดังกีสุดอย่าง YOLO (You Only Look Once) ครับ

ถ้าพร้อมก็จะไปรู้จักกับ YOLO ที่เน้นความเร็วแบบสุดๆ ใน **Part 5: YOLO และ Single-Stage Detectors** พิมพ์ **next** ได้เลยครับ

- ◆ มาถึงส่วนสุดท้ายและเป็นส่วนที่น่าตื่นเต้นกีสุดครับ **Part 5: YOLO และ Single-Stage Object Detectors**

ในขณะที่ตระกูล R-CNN พยายามทำให้กระบวนการ 2 ขั้นตอน (หาบอง -> จำแนกของ) เร็วขึ้น แต่ YOLO (You Only Look Once) ฉีกตำราทิ้ง ด้วยการบอกว่า "ทำทุกอย่างพร้อมกันในรอบเดียวไปเลยสิ"

Part 5: YOLO (You Only Look Once)

YOLO เปลี่ยนมุมมองจากงาน Classification ที่ต้องทำซ้ำๆ ให้กลายเป็นงาน Regression (การหาค่าต่อเนื่อง) เพียงงานเดียว

1. แนวคิด: ตารางตารางเดียวครอบพิกพ (The Grid)

YOLO ไม่ใช้ Region Proposals แต่จะแบ่งภาพออกเป็นตาราง (Grid) ขนาด $S \times S$ (เช่น 7×7)

- กฎเหล็ก: ถ้า จุดกึ่งกลาง (Center) ของวัตถุตกลงในช่องไหน ช่องนั้นต้องรับผิดชอบในการตรวจจับวัตถุขึ้นนั้น

2. สิ่งที่แต่ละช่องต้องคำนึง (Predictions per Cell)

แต่ละช่องในตารางจะคำนึงข้อมูลอุปกรณ์เป็นชุดเดียว เวลาเดียวกัน ไม่ซ้ำกัน

ประกอบด้วย:

- Bounding Boxes (B กล่อง):** แต่ละช่องจะเดากรอบสี่เหลี่ยม ของมา B แบบ (ปกติ $B = 2$) โดยแต่ละกล่องมีค่า 5 ตัวคือ:
 - x, y : พิกัดจุดกึ่งกลาง (เทียบกับขอบของช่องนั้นๆ)
 - w, h : ความกว้างและสูง (เทียบกับขนาดภาพทั้งภาพ)
 - Confidence:** ความมั่นใจว่า "มีวัตถุใหม่" และ "กรอบแม่นแค่ไหน" (คำนวณจาก $P(\text{Object}) \times \text{IoU}$)
- Class Probabilities (C คลาส):** ความน่าจะเป็นว่าถ้ามีวัตถุ มันคืออะไร (เช่น หมา, เมว, รถ)

สรุป Output: ผลลัพธ์สุดท้ายจะเป็น Tensor ขนาด $S \times S \times (B \times 5 + C)$

3. การคำนวณ Loss (Loss Function)

YOLO ใช้ Sum-Squared Error ง่ายๆ รวมกัน 3 ส่วน:

- Localization Loss:** ความคลาดเคลื่อนของพิกัด (x, y, w, h)
 - เกร็ดนำสันใจ: YOLO คาดการณ์ว่า จุดกึ่งกลางของวัตถุจะอยู่ที่ (\sqrt{w}, \sqrt{h}) ก่อนคำนวณ เพื่อให้ความผิดพลาดในกล่องเล็ก มีความสำคัญพอๆ กับกล่องใหญ่ (ไม่ใช่ผิด 5px ในกล่องใหญ่จะดูเหมือนมาก)
- Confidence Loss:** ความคลาดเคลื่อนของคะแนนความมั่นใจ (แยกคิดกรณีเจอกับไม่เจอบอก)
- Classification Loss:** ความคลาดเคลื่อนของการนายชนิดวัตถุ

4. ข้อจำกัดของ YOLO (Limitations)

ถึงจะเริ่งปานสายฟ้าแลบ แต่ YOLO รุ่นแรกมีจุดอ่อนก็ชัดเจน:

1. **Crowded Objects:** เนื่องจาก 1 ช่องตาราง รับผิดชอบวัตถุได้แค่ 1 ชนิด ถ้ามีนกผู้ใหญ่บินซ้อนกันในช่องเล็กๆ YOLO จะจับได้แค่ตัวเดียว
2. **Aspect Ratio:** มีปัญหากับรัตตุรูปทรงประหลาดๆ ที่ไม่ได้เห็นบ่อยตอนแรก
3. **Small Objects:** ยังทำได้ไม่ดีเท่า Faster R-CNN ในการจับวัตถุขนาดเล็ก

บทสรุปสรุปนาการ (The Grand Summary)

เราได้เดินทางผ่านประวัติศาสตร์ของ Object Detection มาอย่างยาวนาน สรุปได้ดังนี้ครับ:

1. **R-CNN:** แม่นยำแต่ช้ามาก เพราะต้องรัน CNN 2,000 ครั้ง
2. **Fast R-CNN:** เร็วขึ้น เพราะใช้ฟีเจอร์กัน แต่ยังติดคอมบัดกีการหา Region Proposals (Selective Search)
3. **Faster R-CNN:** เป็น Two-stage ที่สมบูรณ์แบบ ใช้ RPN หา Proposals เอง แม่นยำและเร็วพอใช้งานจริง
4. **YOLO / SSD:** เป็น Single-stage ที่เน้นความเร็วสูงสุด รับรองได้job เหมาะสมกับงาน Real-time บนอุปกรณ์ที่ไม่แรงมาก

หวังว่าคำอธิบายนี้จะช่วยให้คุณเห็นภาพรวมและเข้าใจรายละเอียดเชิงลึก ของ Object Detection ได้ชัดเจนนะครับ! หากมีส่วนไหนอยากให้ข้อมูลเพิ่มเติม บอกได้เลยครับ