

Scan2CapMMT

Dense Captioning for 3D Scenes
with Transformers

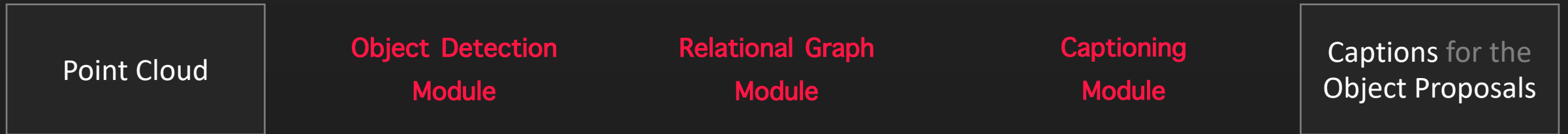
Scan2CapMMT

- I. Scan2CapMMT Recap
- II. Improving Scan2CapMMT
- III. Quantitative & Qualitative Results
- IV. Detection with Transformers
- V. Timeline until the Final Presentation

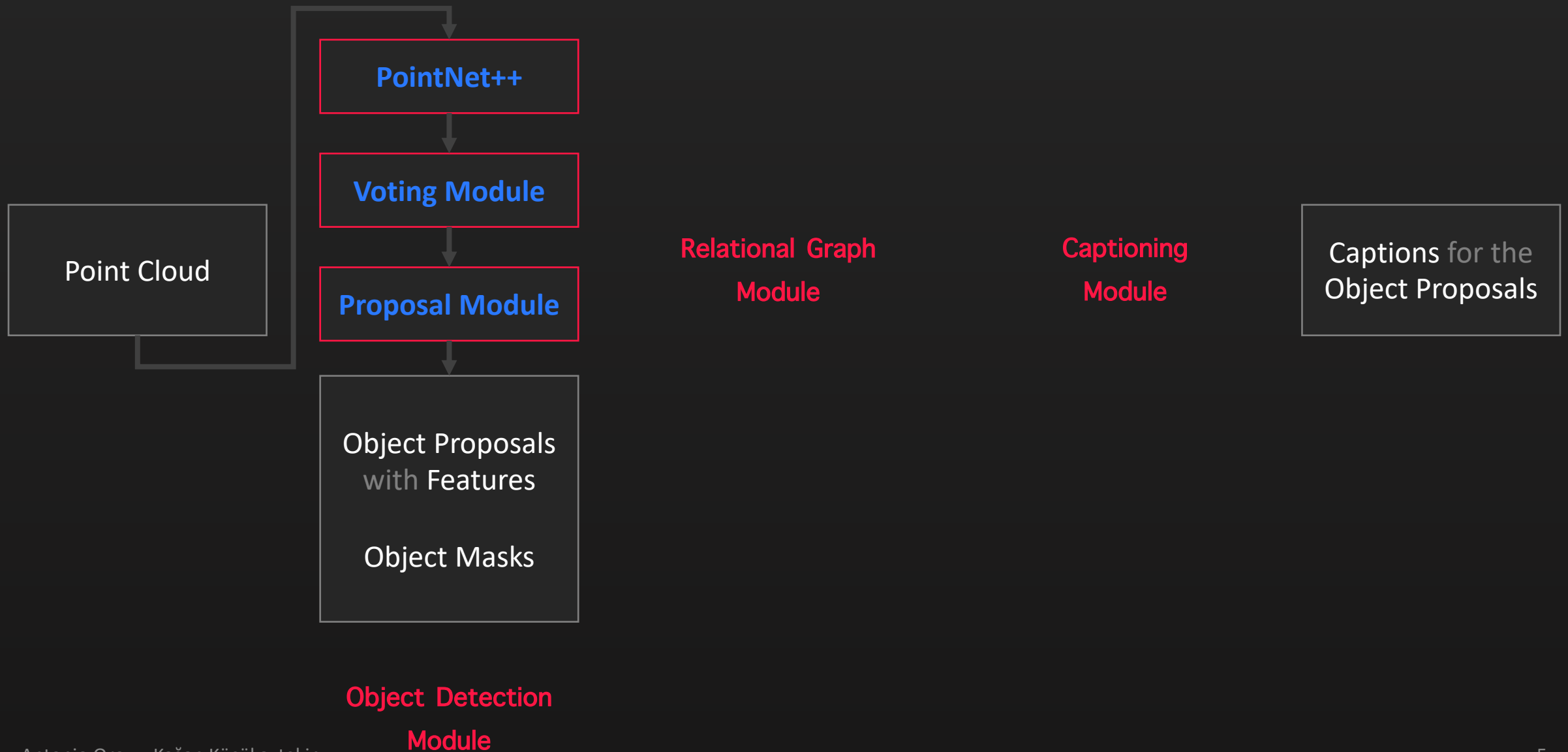
Scan2CapMMT

- I. Scan2CapMMT Recap
- II. Improving Scan2CapMMT
- III. Quantitative & Qualitative Results
- IV. Detection with Transformers
- V. Timeline until the Final Presentation

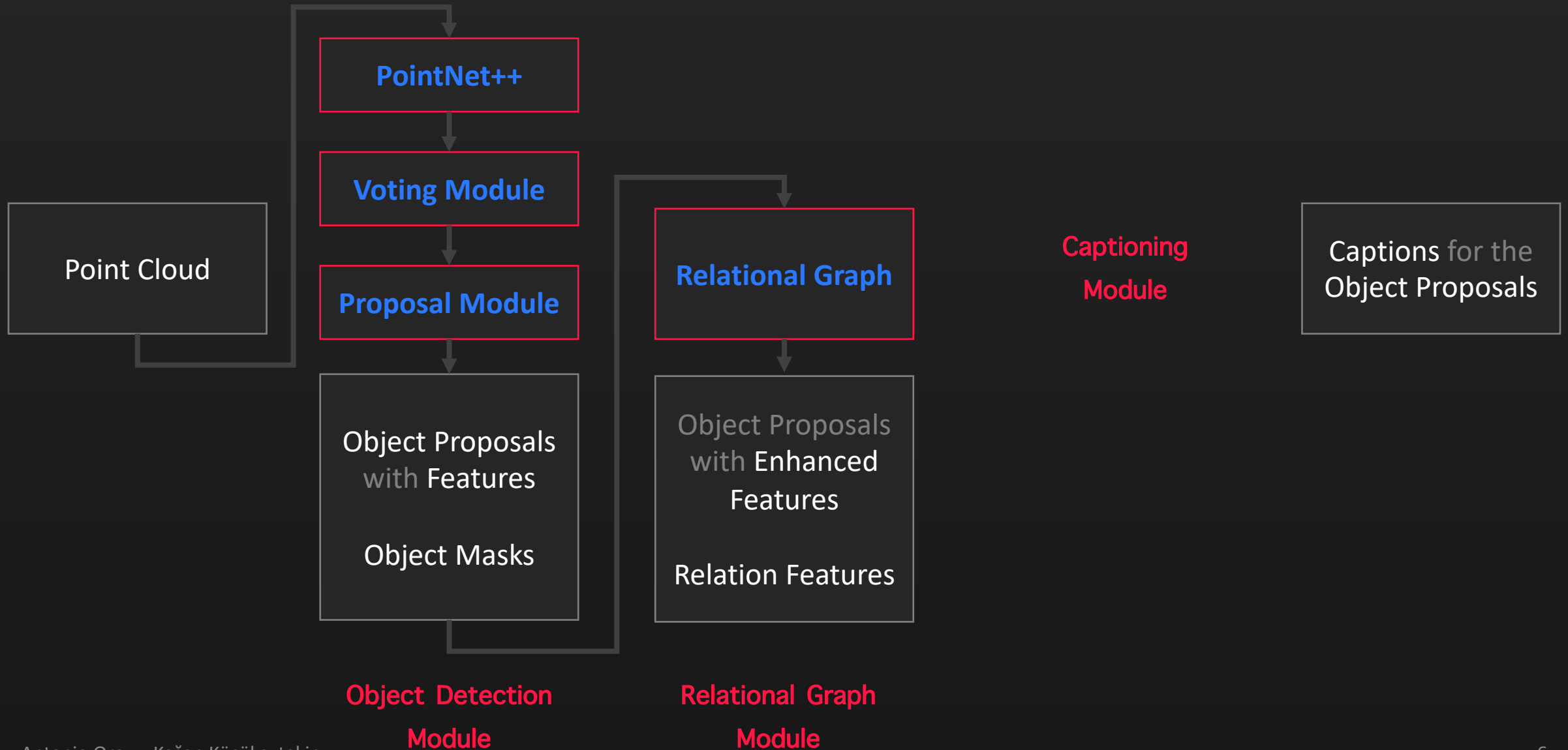
I. Scan2Cap Recap



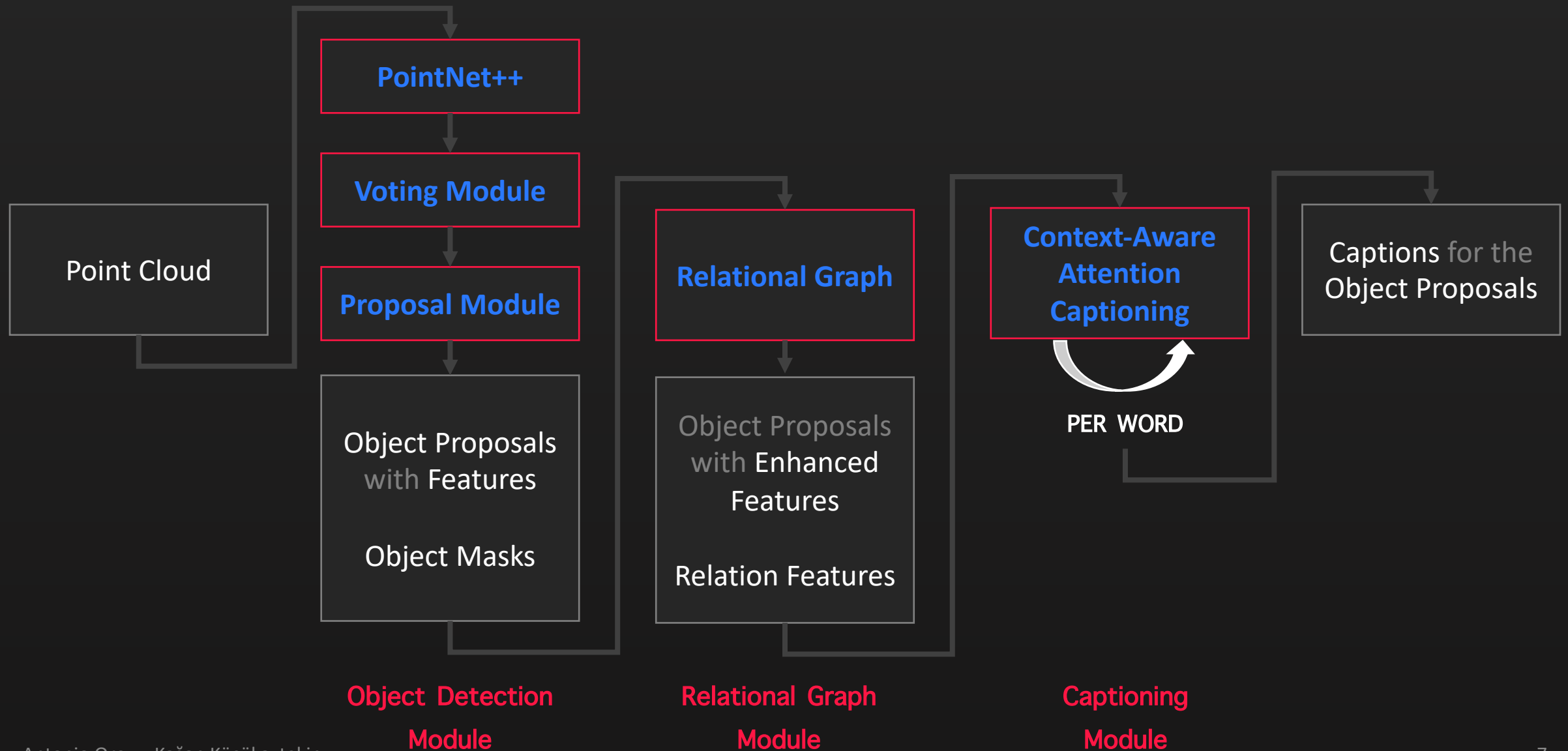
I. Scan2Cap Recap



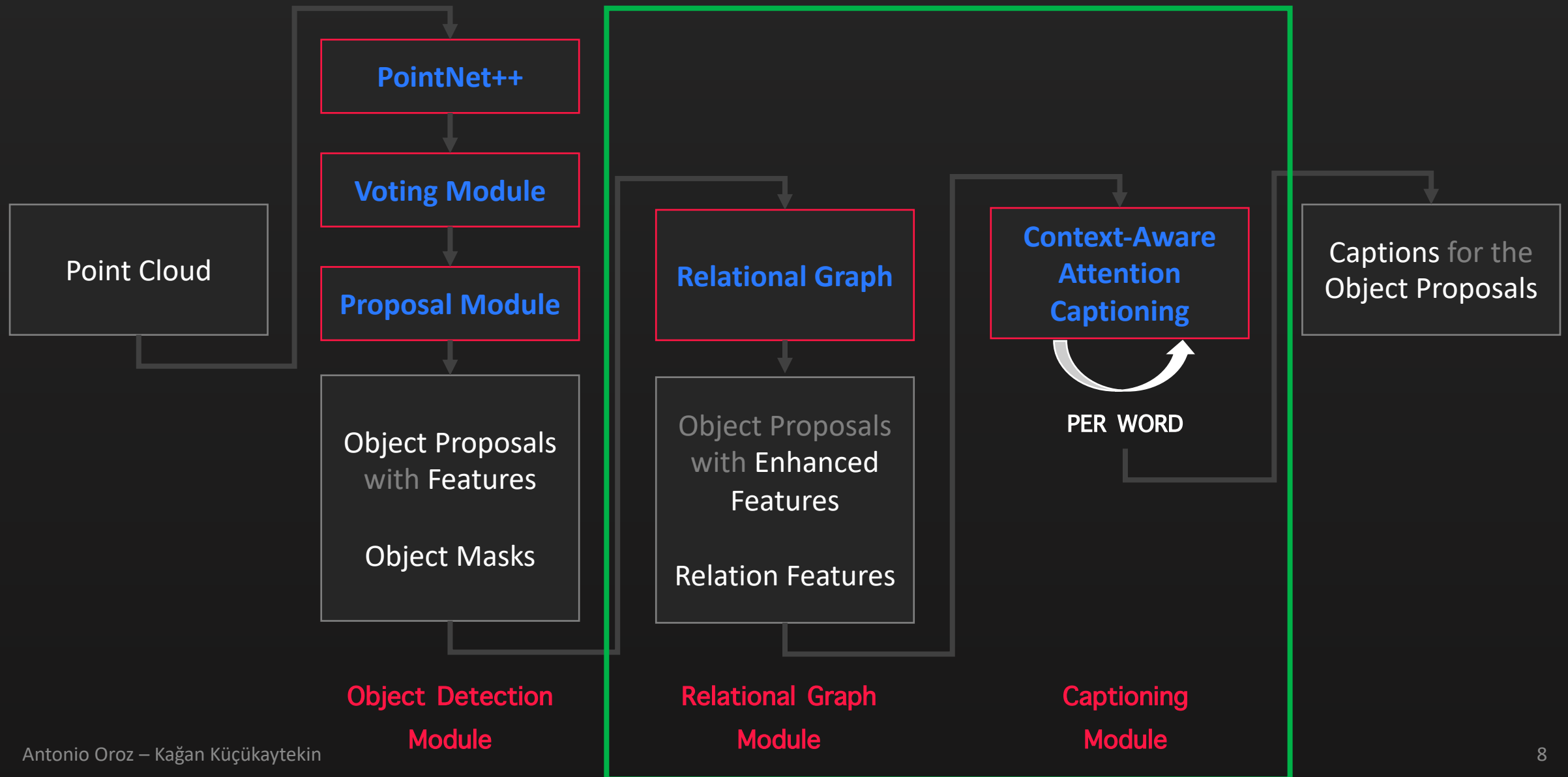
I. Scan2Cap Recap



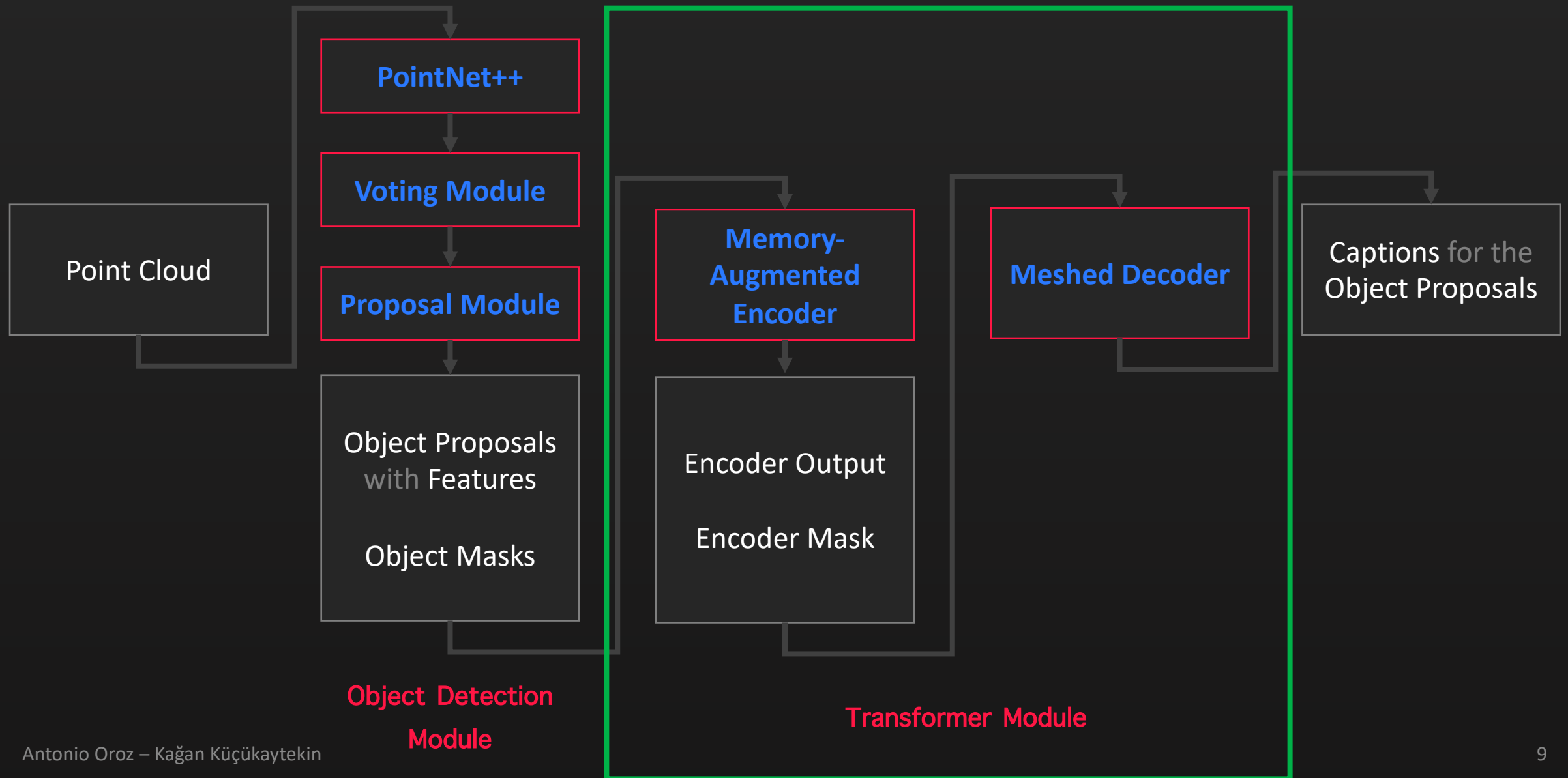
I. Scan2Cap Recap



I. Scan2Cap Recap



I. Scan2CapMMT Recap



Scan2CapMMT

- I. Scan2CapMMT Recap
- II. Improving Scan2CapMMT
- III. Quantitative & Qualitative Results
- IV. Detection with Transformers
- V. Timeline until the Final Presentation

Scan2CapMMT

- I. Scan2CapMMT Recap
- II. Improving Scan2CapMMT
- III. Quantitative & Qualitative Results
- IV. Detection with Transformers
- V. Timeline until the Final Presentation

II. Improving Scan2CapMMT

Beam Search

ITERATIVE SEARCH

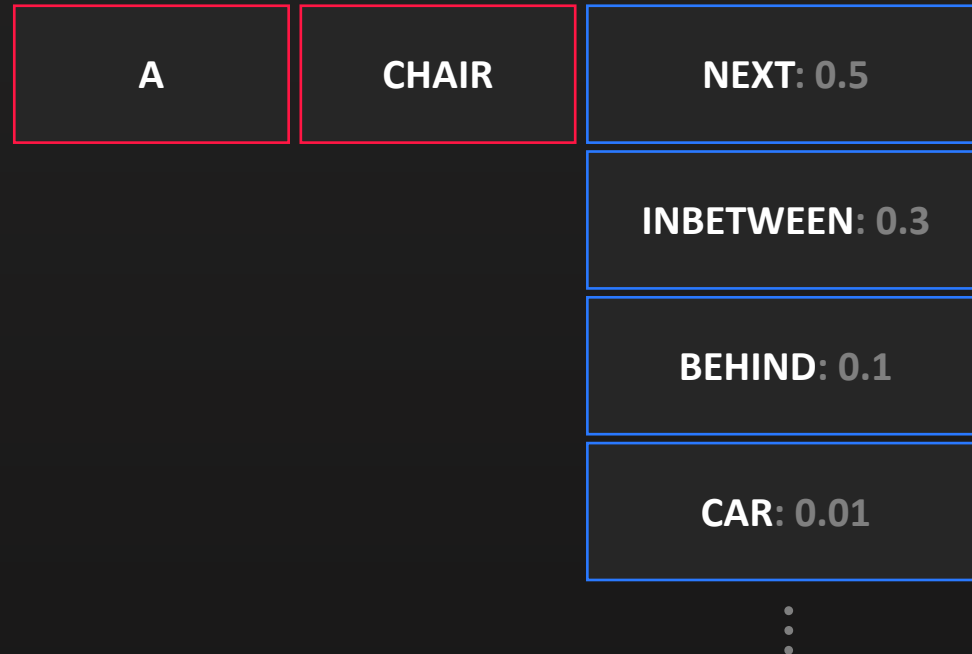
A

CHAIR

II. Improving Scan2CapMMT

Beam Search

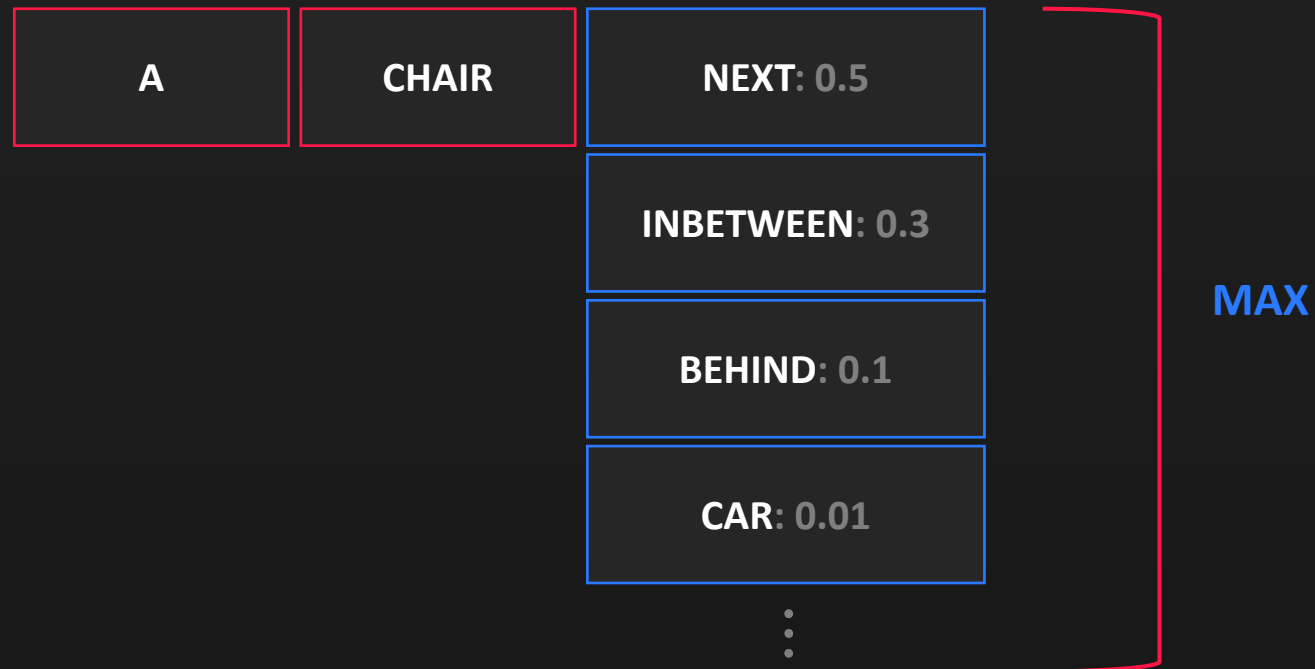
ITERATIVE SEARCH



II. Improving Scan2CapMMT

Beam Search

ITERATIVE SEARCH



II. Improving Scan2CapMMT

Beam Search

ITERATIVE SEARCH

A

CHAIR

NEXT: 0.5

...

II. Improving Scan2CapMMT

Beam Search

A: 0.9

CHAIR: 0.4

BEAM SEARCH
SIZE 2

A: 0.9

TABLE: 0.3

II. Improving Scan2CapMMT

Beam Search

BEAM SEARCH
SIZE 2

A: 0.9

CHAIR: 0.4

NEXT: 0.5

INBETWEEN: 0.3

BEHIND: 0.1

⋮

A: 0.9

TABLE: 0.3

NEXT: 0.3

INBETWEEN: 0.2

⋮

II. Improving Scan2CapMMT

Beam Search

BEAM SEARCH
SIZE 2

A: 0.9	CHAIR: 0.4	NEXT: 0.5	$0.9 * 0.4 * 0.5 = 0.18$
		INBETWEEN: 0.3	$0.9 * 0.4 * 0.3 = 0.108$
		BEHIND: 0.1	$0.9 * 0.4 * 0.1 = 0.036$
		⋮	
A: 0.9	TABLE: 0.3	NEXT: 0.3	$0.9 * 0.3 * 0.3 = 0.081$
		INBETWEEN: 0.2	$0.9 * 0.3 * 0.2 = 0.054$
		⋮	

II. Improving Scan2CapMMT

Beam Search

BEAM SEARCH
SIZE 2



II. Improving Scan2CapMMT

Beam Search

BEAM SEARCH
SIZE 2



II. Improving Scan2CapMMT

Beam Search

A: 0.9

CHAIR: 0.4

NEXT: 0.5

...

BEAM SEARCH
SIZE 2

A: 0.9

CHAIR: 0.4

INBETWEEN: 0.3

...

II. Improving Scan2CapMMT

Beam Search

Reinforcement Learning

CIDEr

This is a white sink... 0.41

This white a rectangular... 0.32

This is kitchen white... 0.24

This white a to sink... 0.001

This is is white oven... 0.09

II. Improving Scan2CapMMT

Beam Search

Reinforcement Learning

CIDEr

This is a white sink... 0.41

This white a rectangular... 0.32

This is kitchen white... 0.24

This white a to sink... 0.001

This is is white oven... 0.09

$$-\frac{1}{k} \sum_{i=1}^k (r(w^i) - b) \log(p(w^i))$$

II. Improving Scan2CapMMT

Beam Search

Reinforcement Learning

CIDEr

This is a white sink... 0.41

This white a rectangular... 0.32

This is kitchen white... 0.24

This white a to sink... 0.001

This is is white oven... 0.09

MEAN
 $b=0.21$

$$-\frac{1}{k} \sum_{i=1}^k (r(w^i) - b) \log(p(w^i))$$

II. Improving Scan2CapMMT

Beam Search

Reinforcement Learning

$$r(w^i) - b$$

This is a white sink... 0.2

This white a rectangular... 0.11

This is kitchen white... 0.03

This white a to sink... -0.21

This is is white oven... -0.12

MEAN
 $b=0.21$

$$-\frac{1}{k} \sum_{i=1}^k (r(w^i) - b) \log(p(w^i))$$

Scan2CapMMT

- I. Scan2CapMMT Recap
- II. Improving Scan2CapMMT
- III. Quantitative & Qualitative Results
- IV. Detection with Transformers
- V. Timeline until the Final Presentation

Scan2CapMMT

- I. Scan2CapMMT Recap
- II. Improving Scan2CapMMT
- III. Quantitative & Qualitative Results
- IV. Detection with Transformers
- V. Timeline until the Final Presentation

III. Quantitative Results: vs Scan2Cap

@0.5IoU

Model	CIDEr	Bleu-4	Meteor	Rouge
VoteNet+GRU	34.31	21.42	20.13	41.33
VoteNet+MMT <i>=Scan2CapMMT</i>	32.99	21.92	20.96	44.40

III. Quantitative Results: vs Scan2Cap

@0.5IoU

Model	CIDEr	Bleu-4	Meteor	Rouge
VoteNet+CAC	36.15	21.58	20.65	41.78
VoteNet+MMT <i>=Scan2CapMMT</i>	32.99	21.92	20.96	44.40

III. Quantitative Results: vs Scan2Cap

@0.5IoU

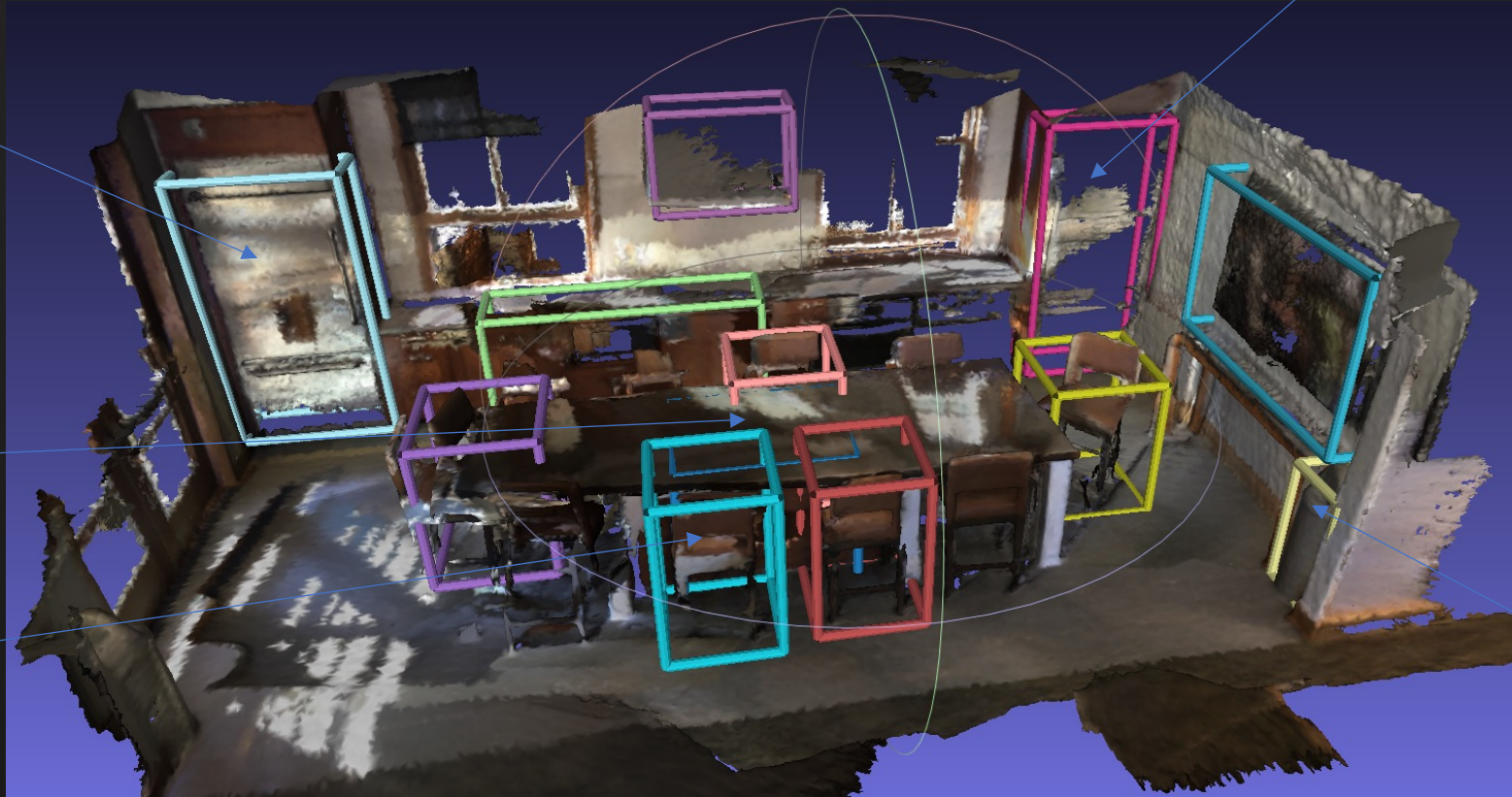
Model	CIDEr	Bleu-4	Meteor	Rouge
VoteNet+RG+CAC <i>= Scan2Cap</i>	39.08	23.32	21.97	44.78
Scan2CapMMT <i>= Scan2CapMMT</i>	32.99	21.92	20.96	44.40

III. Quantitative Results: Reinforcement Learning

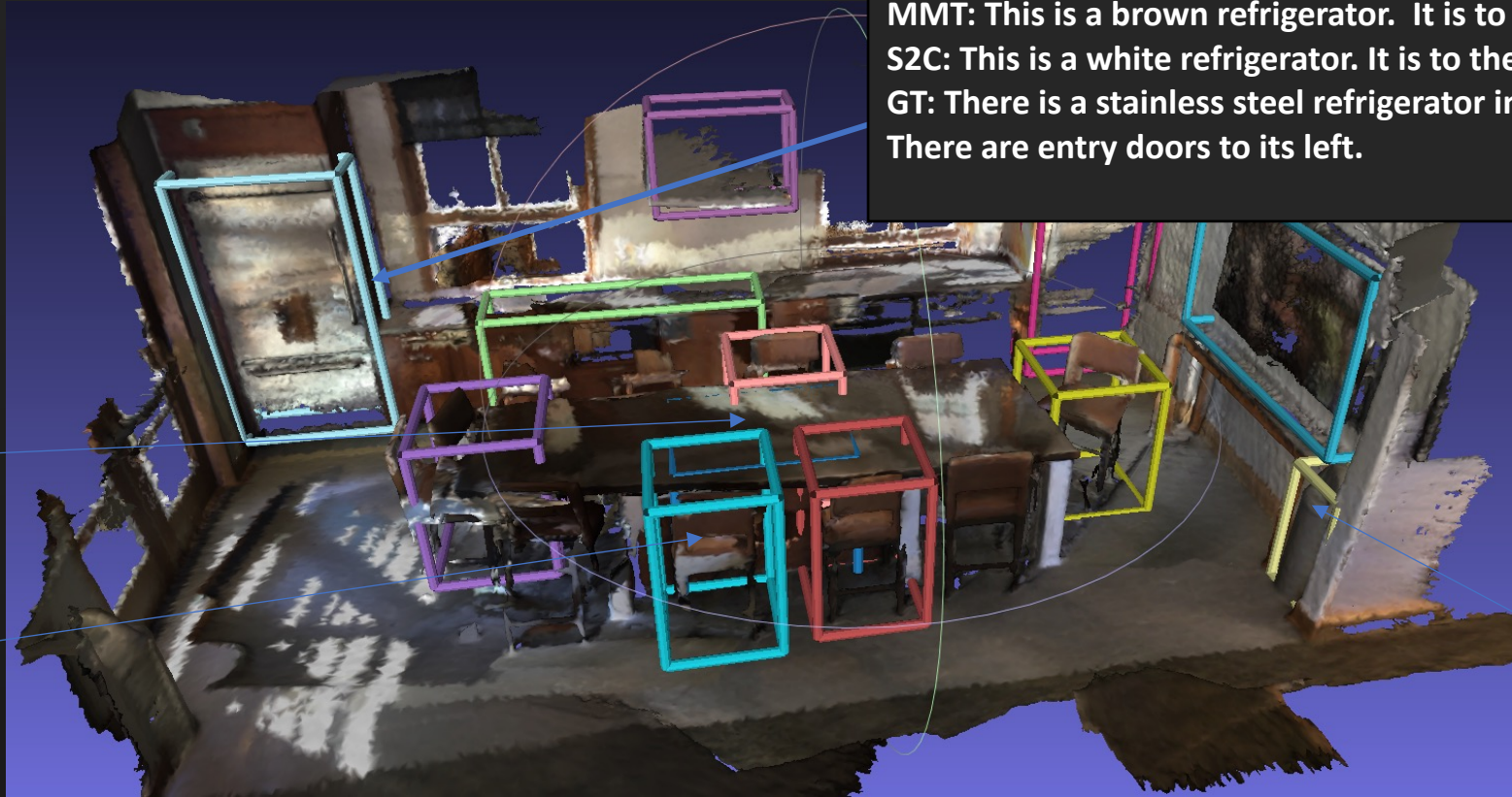
@0.5IoU

Model	CIDEr	Bleu-4	Meteor	Rouge
VoteNet+RG+CAC <i>= Scan2Cap</i>	39.08	23.32	21.97	44.78
VoteNet+MMT <i>= Scan2CapMMT</i>	32.99	21.92	20.96	44.40
Scan2CapMMT RL	36.18	23.68	21.33	44.64

III. Qualitative Results



III. Qualitative Results



MMT: This is a brown refrigerator. It is to the left of the door.
S2C: This is a white refrigerator. It is to the left of the door.
GT: There is a stainless steel refrigerator in corner of the room.
There are entry doors to its left.

III. Qualitative Results

MMT: This is a brown refrigerator. It is to the left of the door.

S2C: This is a white refrigerator. It is to the left of the door.

GT: There is a stainless steel refrigerator in corner of the room. There are entry doors to its left.

MMT: This is a black chair. It is at the table.

S2C: This is a wooden chair. It is to the left of another chair.

GT: This is a brown chair. It is in between two other chairs.



III. Qualitative Results

MMT: This is a brown refrigerator. It is to the left of the door.

S2C: This is a white refrigerator. It is to the left of the door.

GT: There is a stainless steel refrigerator in corner of the room. There are entry doors to its left.

MMT: This is a brown table. It is in front of a window

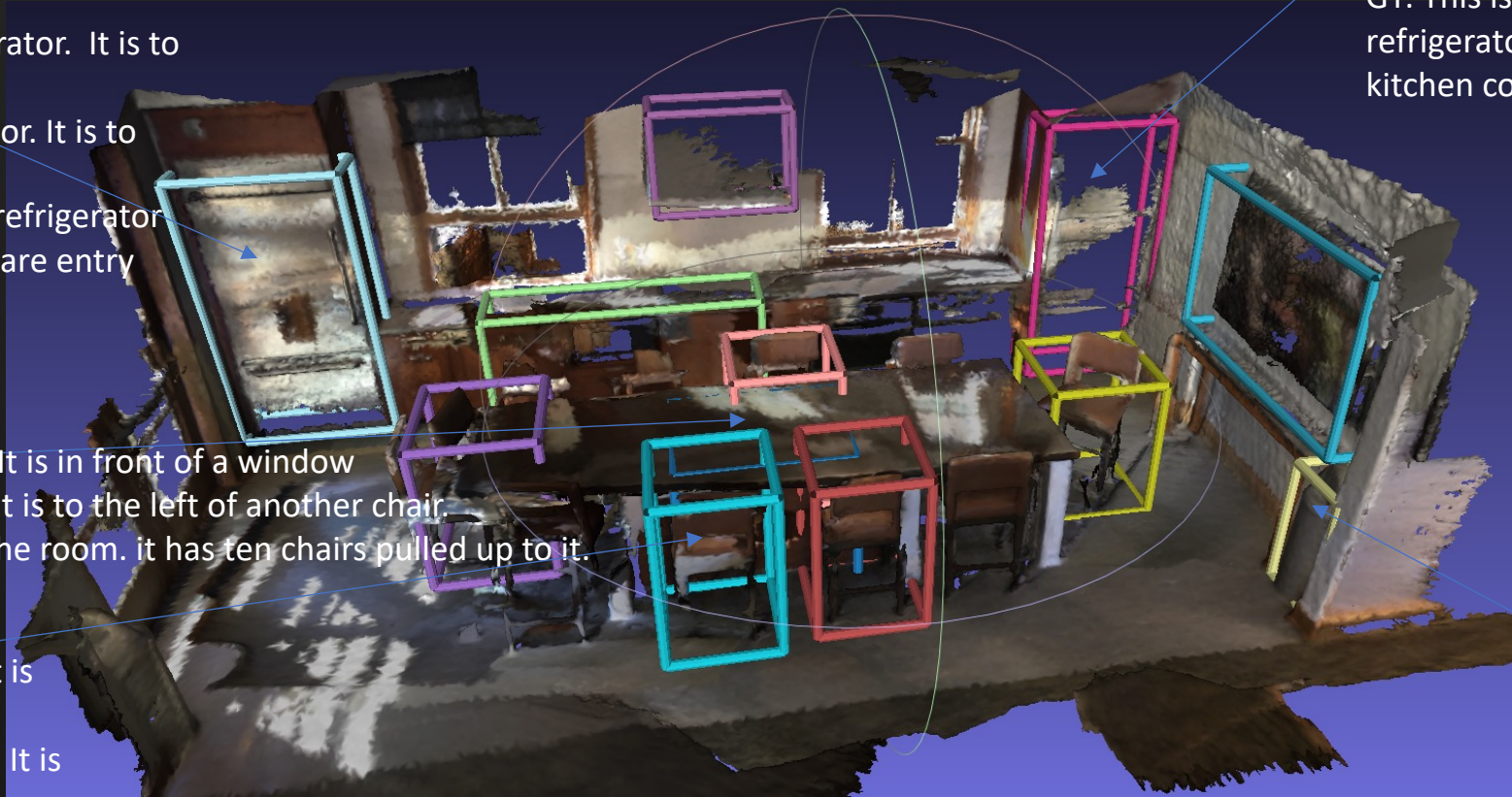
S2C: This is a wooden chair. It is to the left of another chair.

GT: there is a large table in the room. it has ten chairs pulled up to it.

MMT: This is a black chair. It is at the table.

S2C: This is a wooden chair. It is to the left of another chair.

GT: This is a brown chair. It is in between two other chairs.



MMT: This is a brown door. It is to the right of the door.

S2C: This is a white door. It is to the left of the shelf.

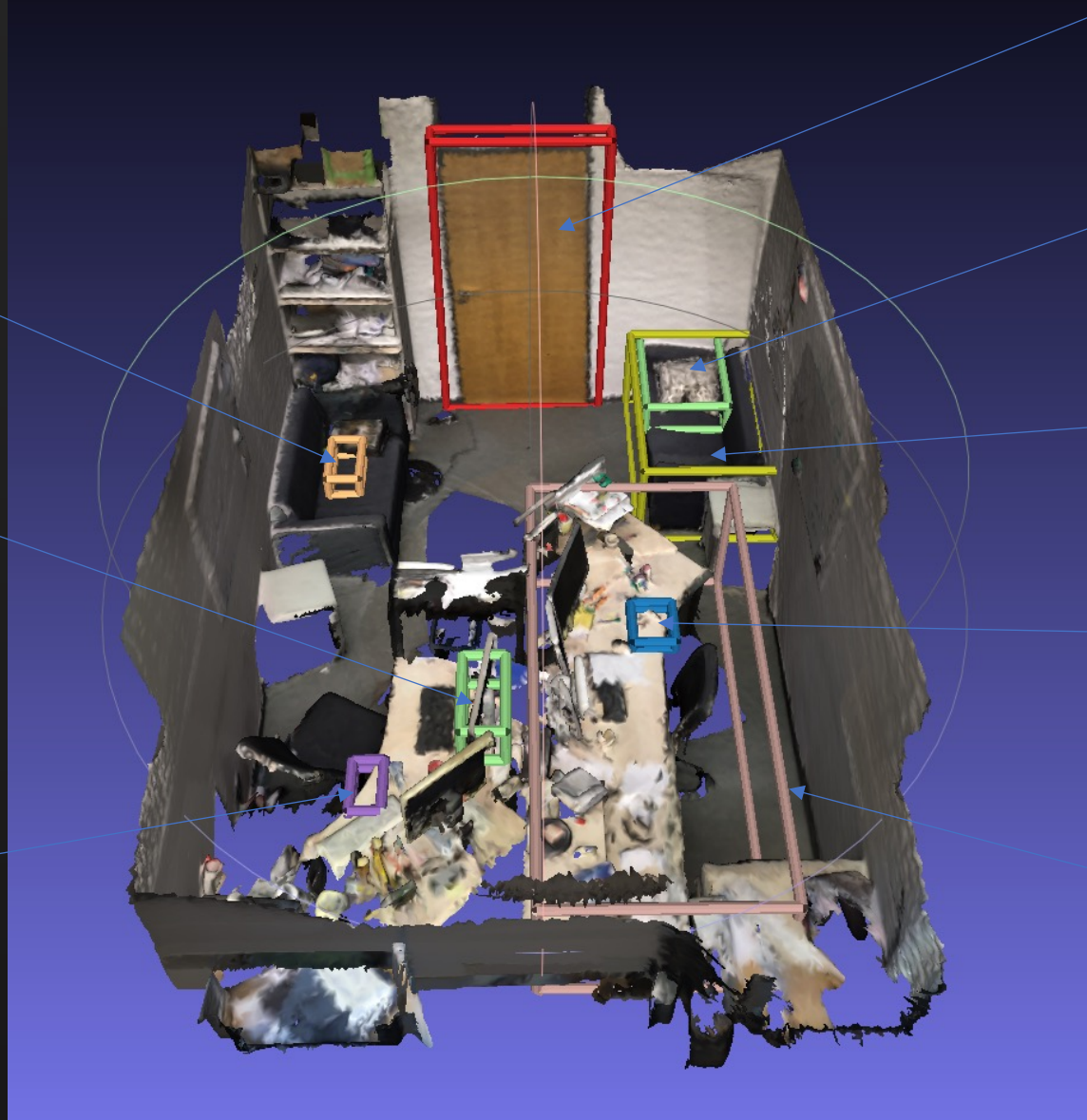
GT: This is a stainless steel refrigerator. It is to the right of a kitchen counter.

MMT: This is a black trash can. It is to the right of the door.

S2C: This is a trash can. It sets against the wall.

GT: This is a gray trash can. It is to the right of a table.

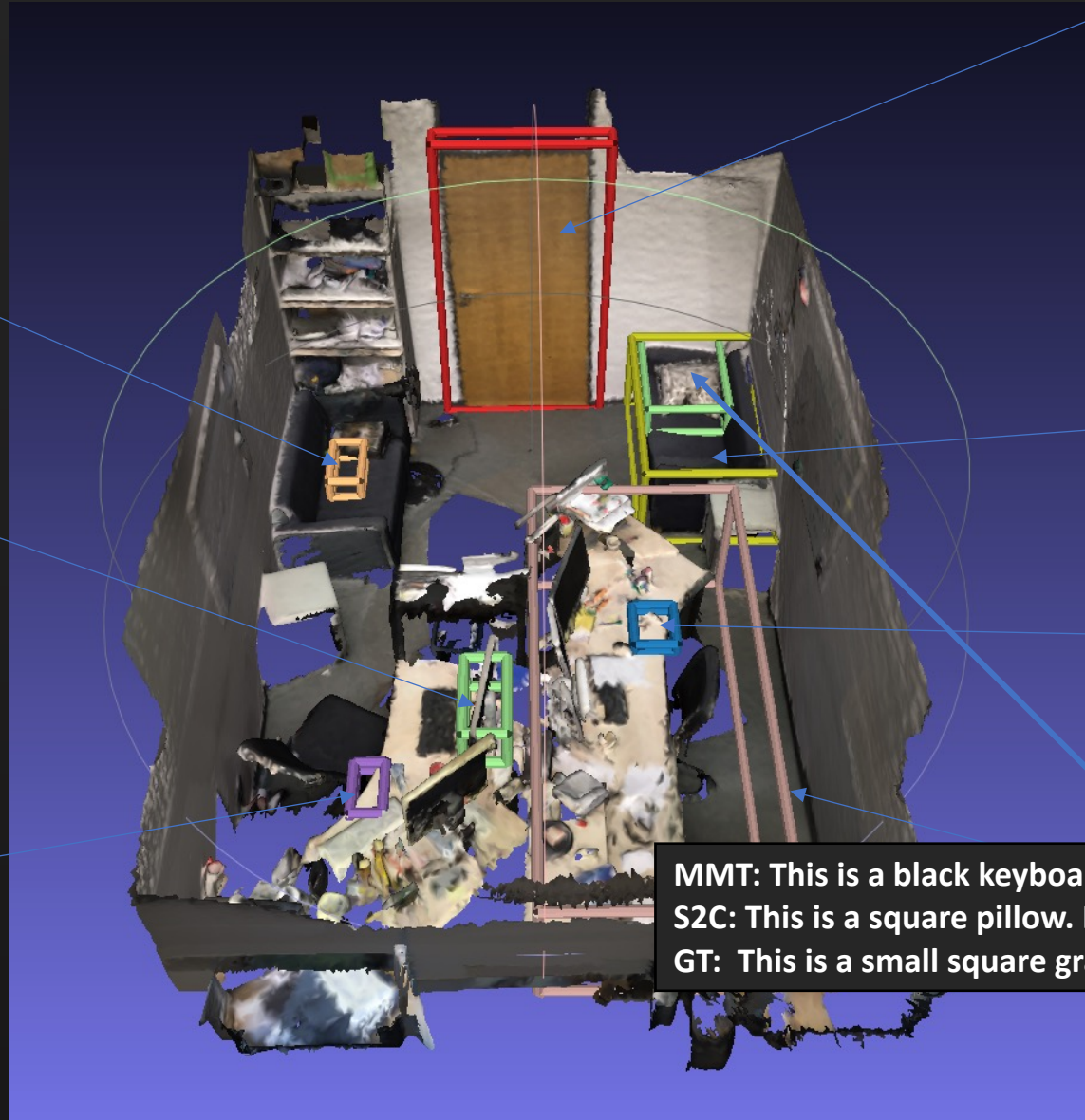
III. Qualitative Results



III. Qualitative Results



III. Qualitative Results



MMT: This is a brown door. It is to the right of a bookshelf.

S2C: This is a white door. It is to the left of a couch.

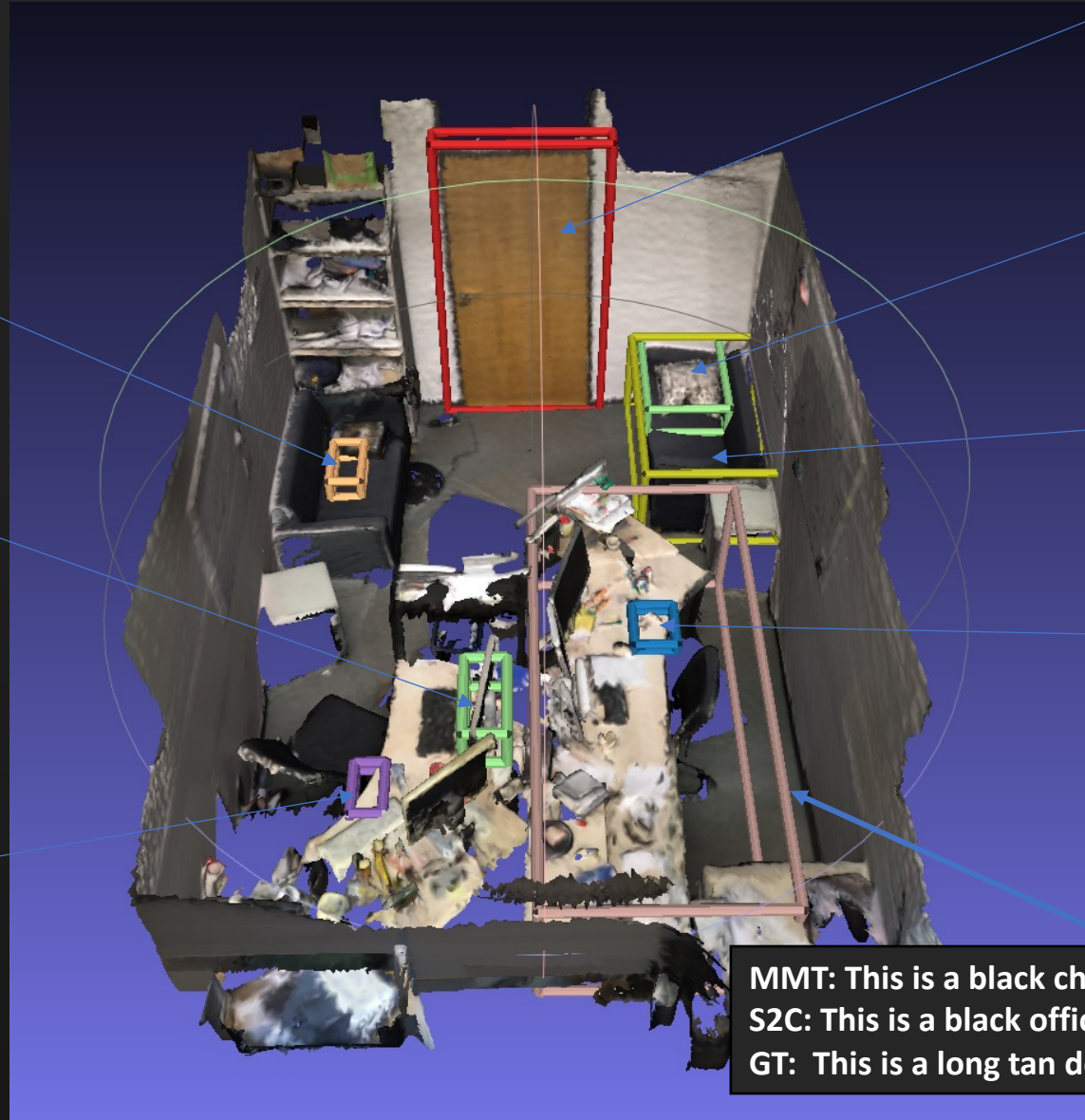
GT: A light brown door beside a tall shelf. A black couch is to the right of it .

MMT: This is a black keyboard. It is on the desk.

S2C: This is a square pillow. It is on the couch.

GT: This is a small square gray pillow. It is located on a black couch.

III. Qualitative Results



MMT: This is a brown door. It is to the right of a bookshelf.

S2C: This is a white door. It is to the left of a couch.

GT: A light brown door beside a tall shelf. A black couch is to the right of it .

MMT: This is a black keyboard. It is on the desk.

S2C: This is a square pillow. It is on the couch.

GT: This is a small square gray pillow. It is located on a black couch.

MMT: This is a black chair. It is to the right of the desk.

S2C: This is a black office chair. It is in front of a desk.

GT: This is a long tan desk. It is located near a wall and a small cabinet.

III. Qualitative Results

MMT: This is a black chair. It is to the right of the desk.

S2C: This is a brown couch. It is to the left of a brown table.

GT: It is a black sofa. It is located to the wall behind the fan.

MMT: This is a black monitor. It is on the desk.

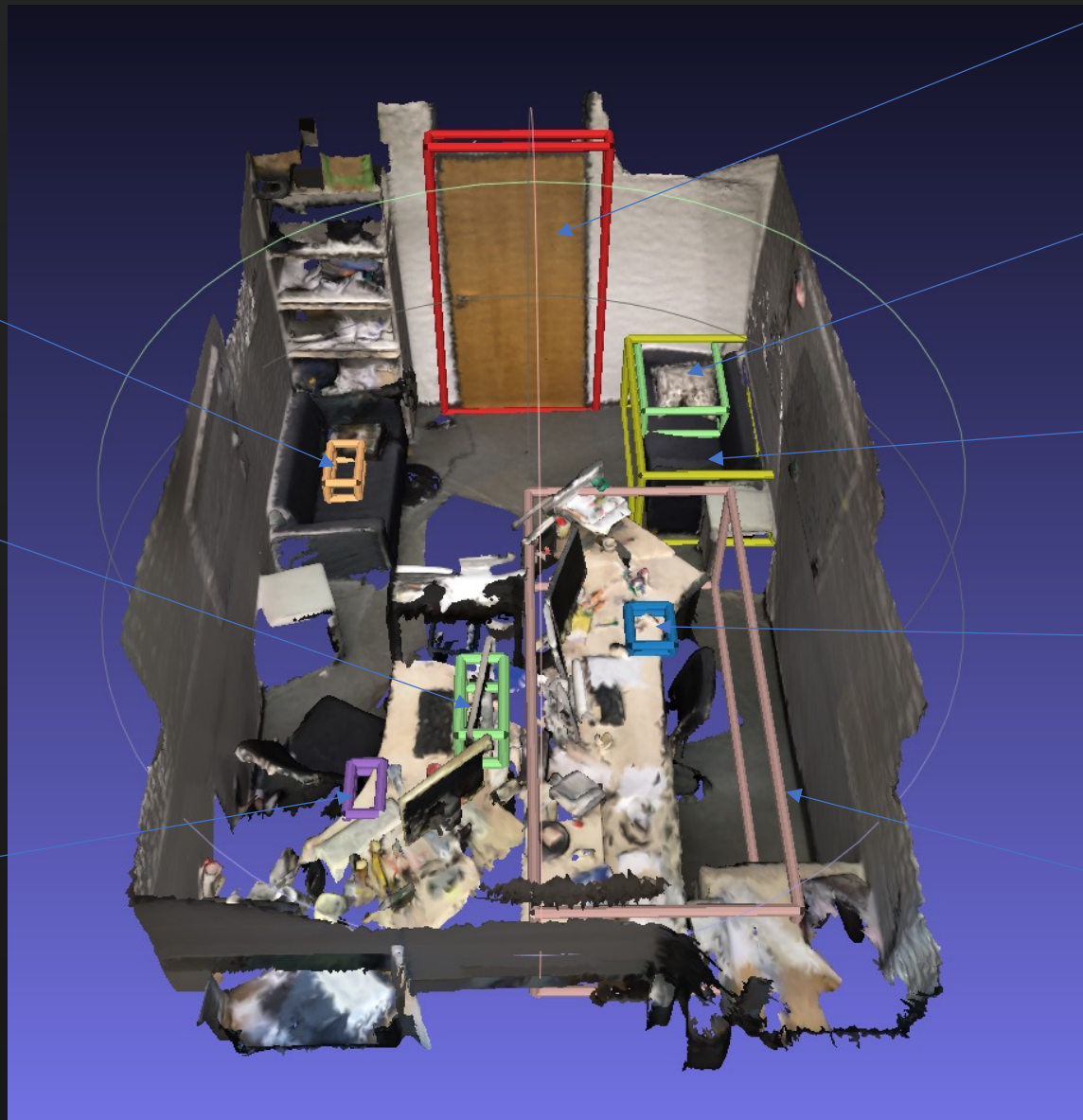
S2C: N/A

GT: The monitor is located on top of the desk, and to the left of the other monitor facing the chair.

MMT: This is a black keyboard. It is on the desk.

S2C: N/A

GT: This is a long tan desk. It is located next to a black office chair.



MMT: This is a brown door. It is to the right of a bookshelf.

S2C: This is a white door. It is to the left of a couch.

GT: A light brown door beside a tall shelf. A black couch is to the right of it .

MMT: This is a black keyboard. It is on the desk.

S2C: This is a square pillow. It is on the couch.

GT: This is a small square gray pillow. It is located on a black couch.

MMT: This is a brown couch. It is to the right of the desk.

S2C: This is a brown couch. It is to the left of a table.

GT: The couch is located in the corner of the room. It is to the right side of the door.

2.MMT: This is a black keyboard. It is on a desk.

S2C: This is a black monitor. It is on a desk.

GT: A black computer screen is sitting on the desk. It is next to a black framed computer screen and to the left of it.

MMT: This is a black chair. It is to the right of the desk.

S2C: This is a black office chair. It is in front of a desk.

GT: This is a long tan desk. It is located near a wall and a small cabinet.

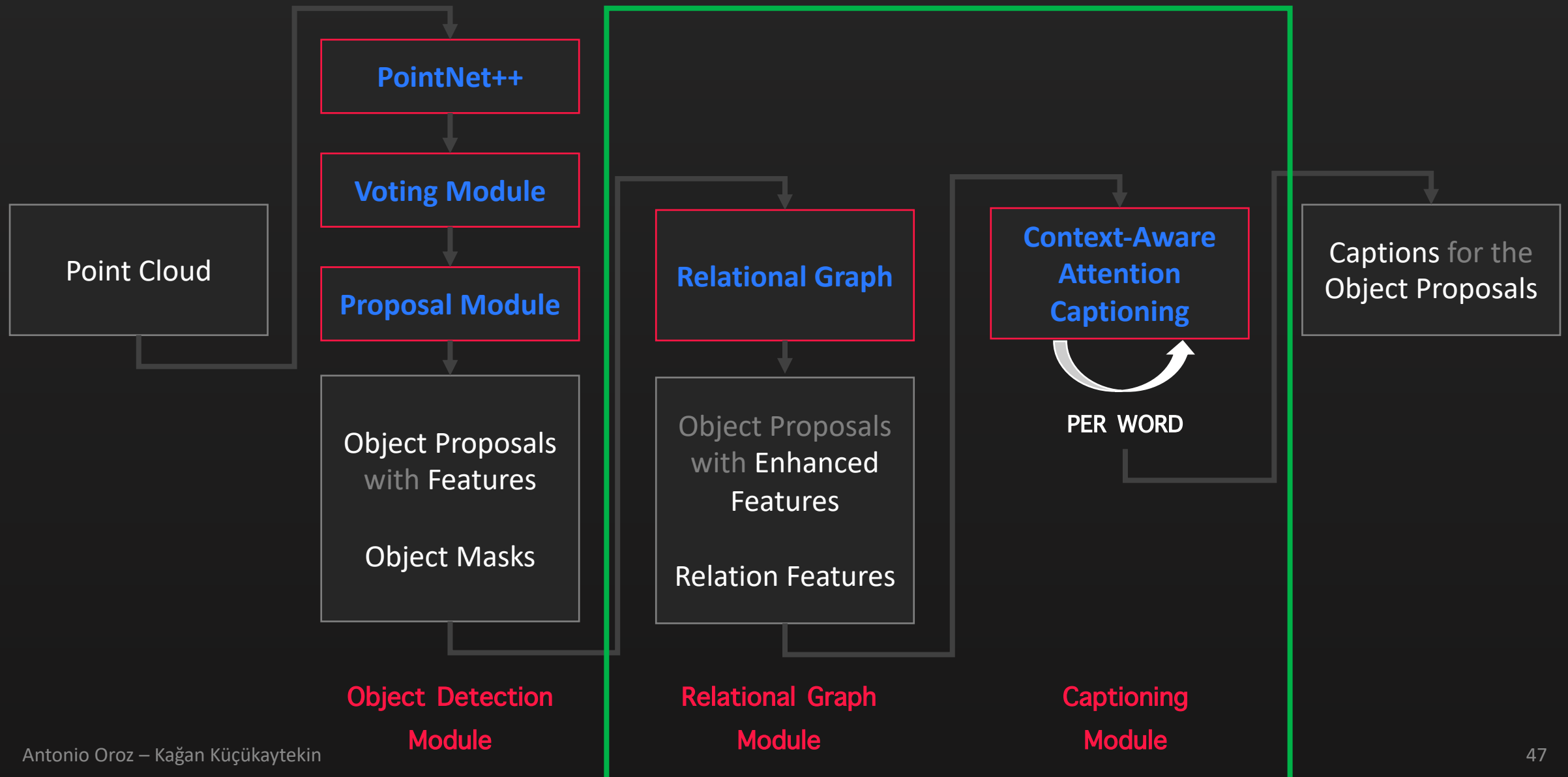
Scan2CapMMT

- I. Scan2CapMMT Recap
- II. Improving Scan2CapMMT
- III. Quantitative & Qualitative Results
- IV. Detection with Transformers
- V. Timeline until the Final Presentation

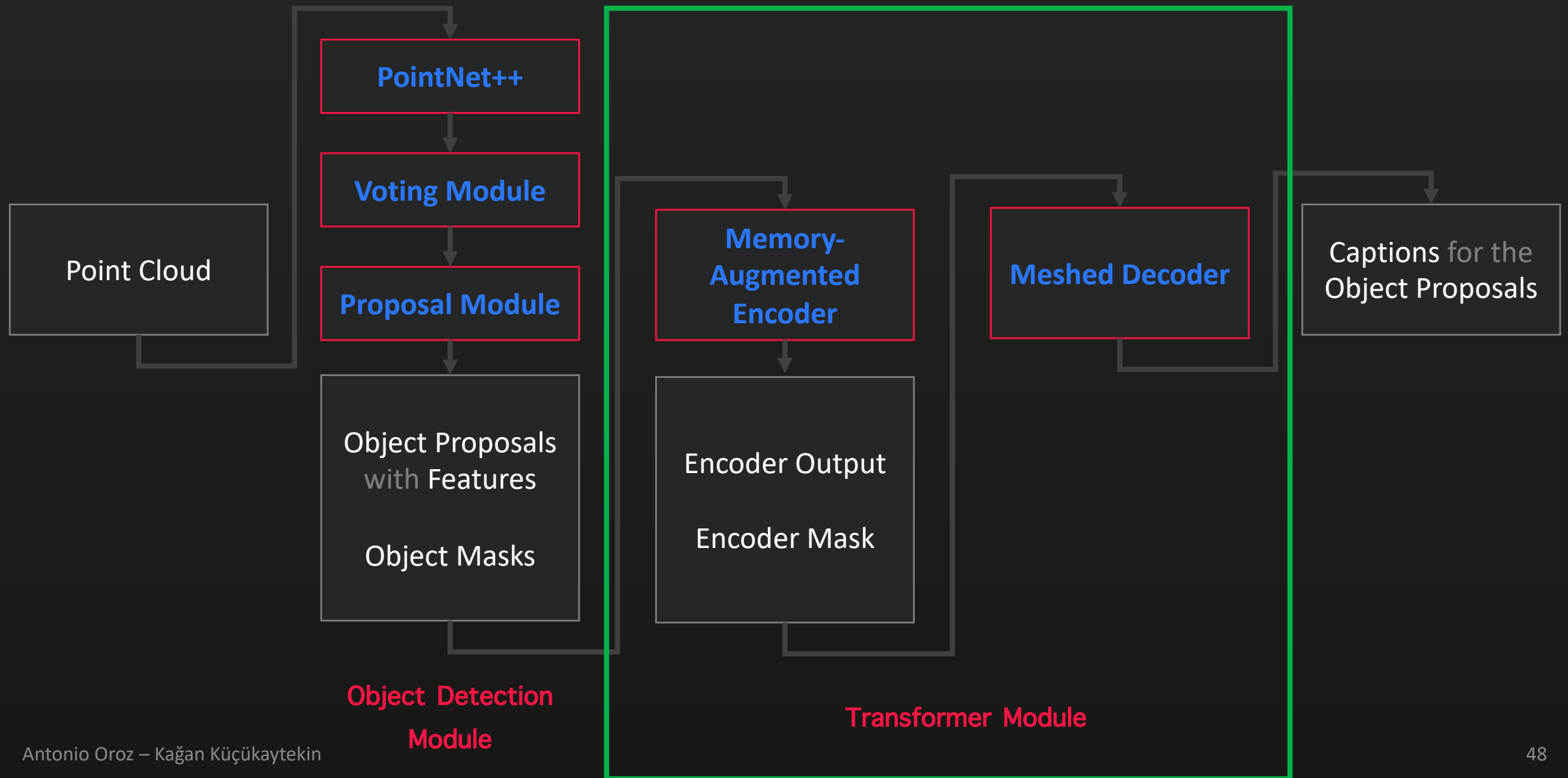
Scan2CapMMT

- I. Scan2CapMMT Recap
- II. Improving Scan2CapMMT
- III. Quantitative & Qualitative Results
- IV. Detection with Transformers
- V. Timeline until the Final Presentation

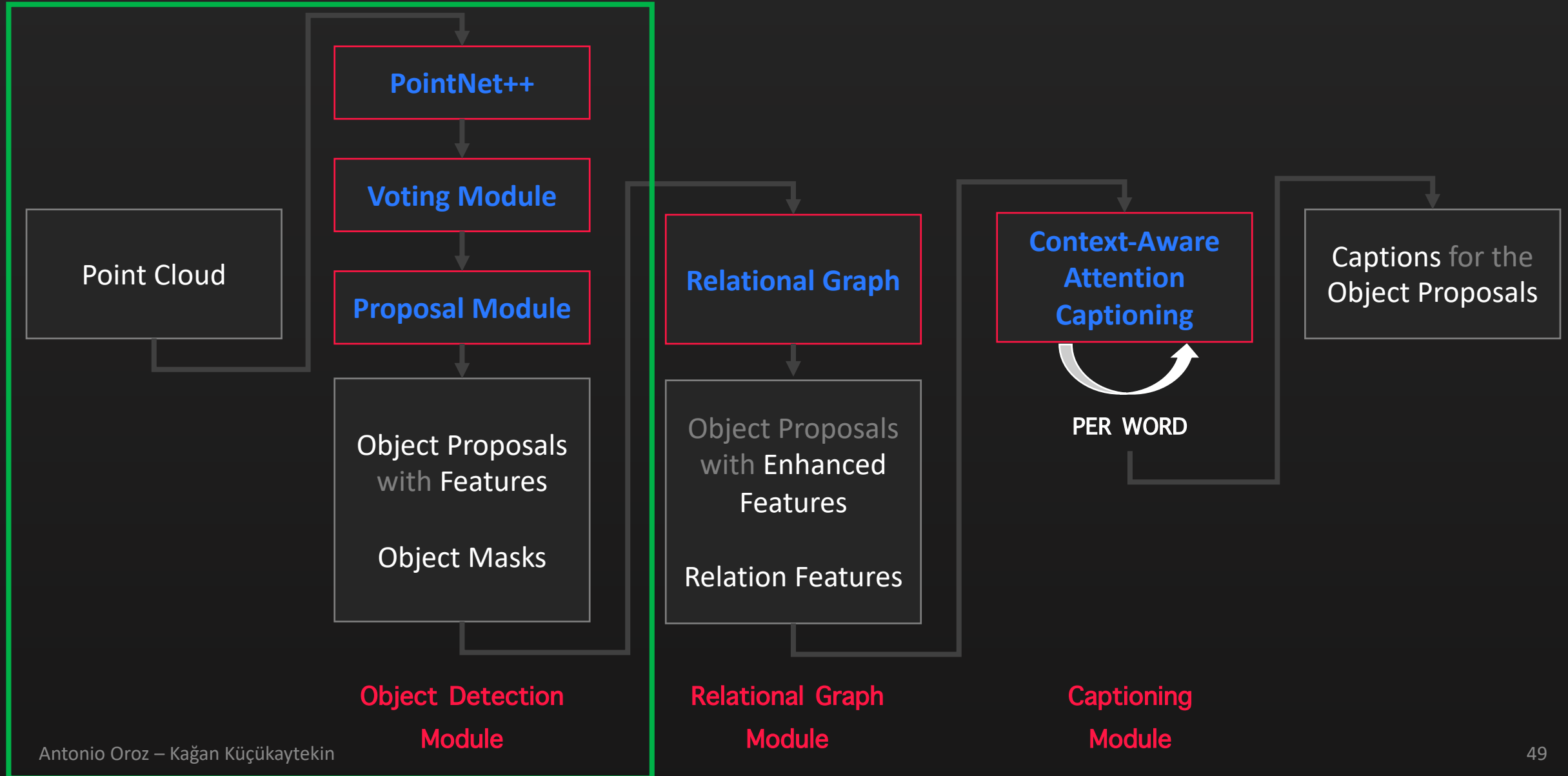
IV. Our Contribution



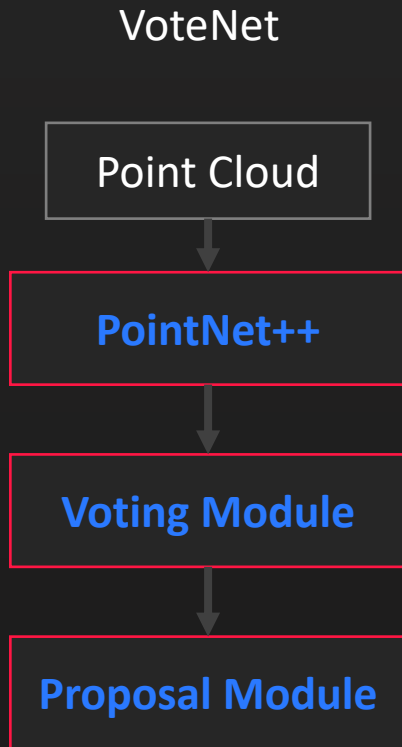
IV. Our Contribution



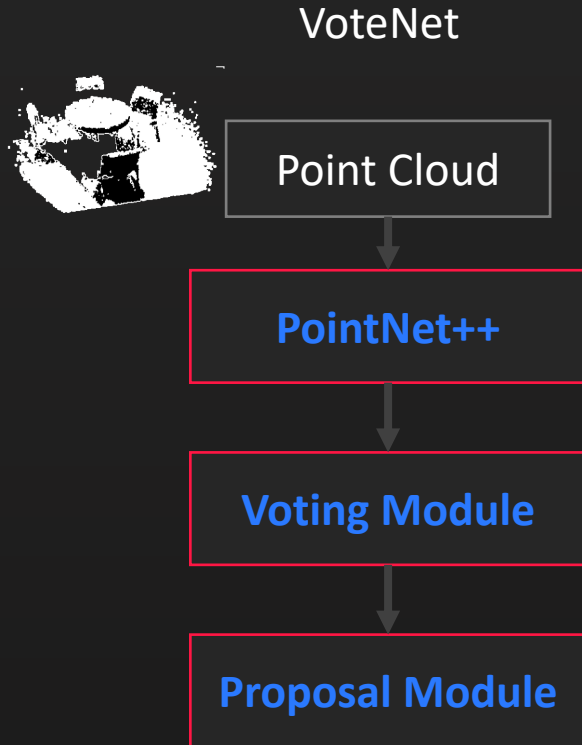
IV. Exploring Transformers for Detection Module



IV. Detection with Transformers



IV. Detection with Transformers



IV. Detection with Transformers

VoteNet

Point Cloud

PointNet++

Voting Module

Proposal Module

IV. Detection with Transformers

VoteNet

Point Cloud

PointNet++

Voting Module

Proposal Module

IV. Detection with Transformers

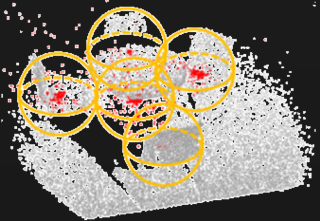
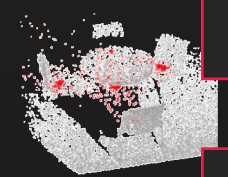
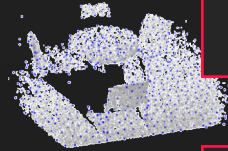
VoteNet

Point Cloud

PointNet++

Voting Module

Proposal Module



IV. Detection with Transformers

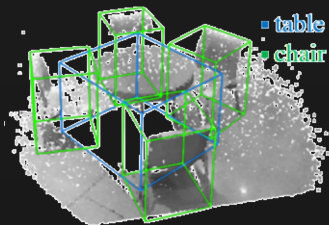
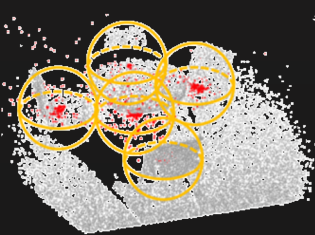
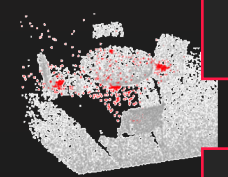
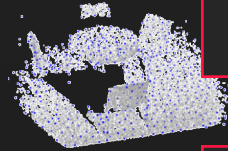
VoteNet

Point Cloud

PointNet++

Voting Module

Proposal Module



IV. Detection with Transformers

VoteNet

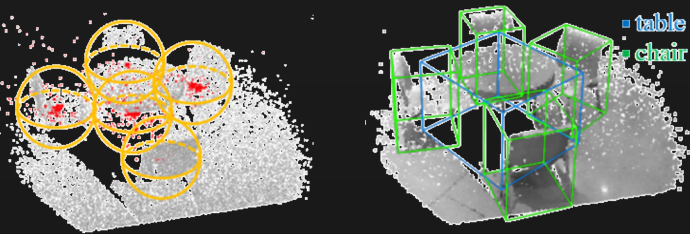
Point Cloud

PointNet++

Voting Module

Proposal Module

*Vote grouping is an issue! Especially when objects are overlapping.



IV. Detection with Transformers

VoteNet

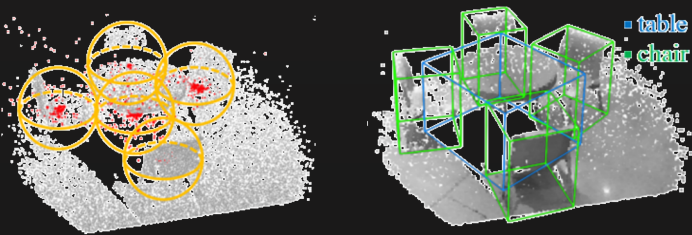
Point Cloud

PointNet++

Voting Module

Proposal Module

- *Vote grouping is an issue! Especially when objects are overlapping.
- *Radius for grouping is an important hyperparameter



IV. Detection with Transformers

VoteNet

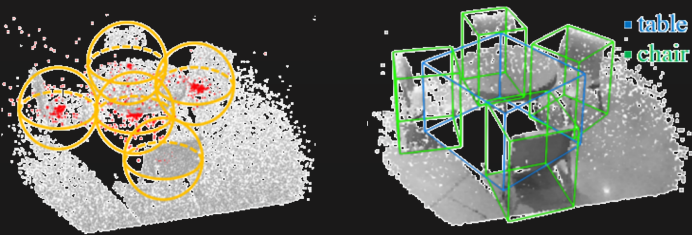
Point Cloud

PointNet++

Voting Module

Proposal Module

- *Vote grouping is an issue! Especially when objects are overlapping.
- *Radius for grouping is an important hyperparameter
- *NMS only for eval



IV. Detection with Transformers

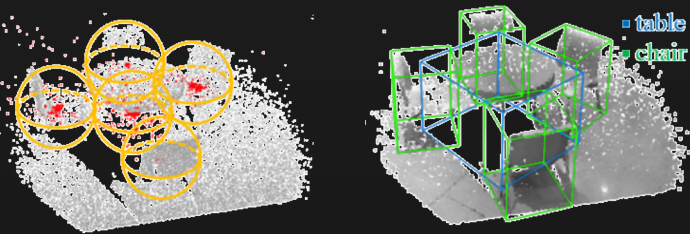
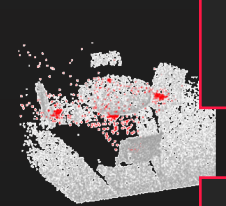
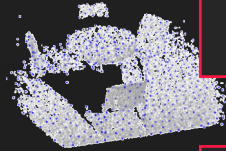
VoteNet

Point Cloud

PointNet++

Voting

Proposal



3DETR

Point Cloud



Group-Free-3D

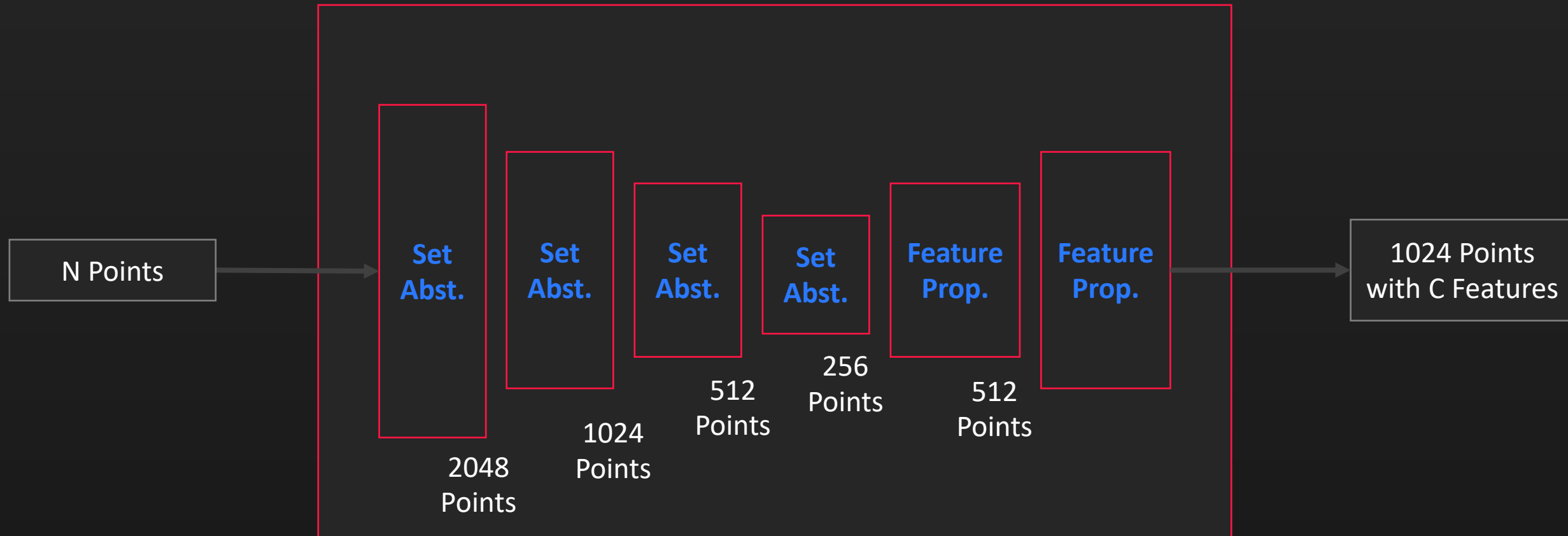
Point Cloud



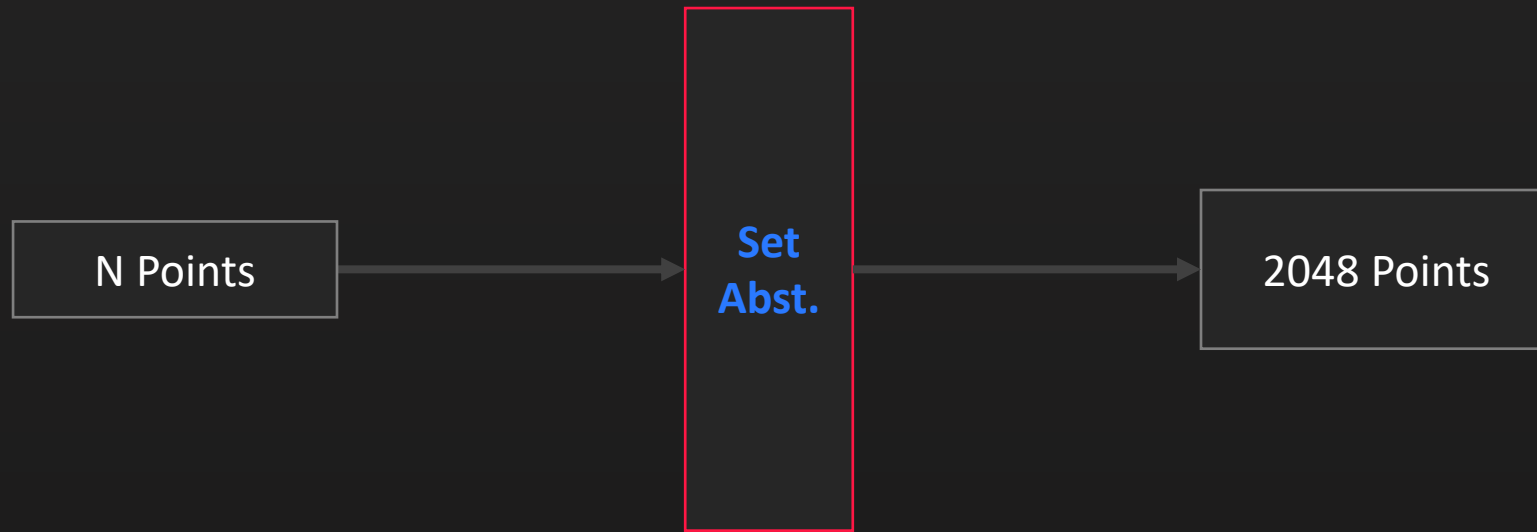
IV. Detection with Transformers

PointNet++

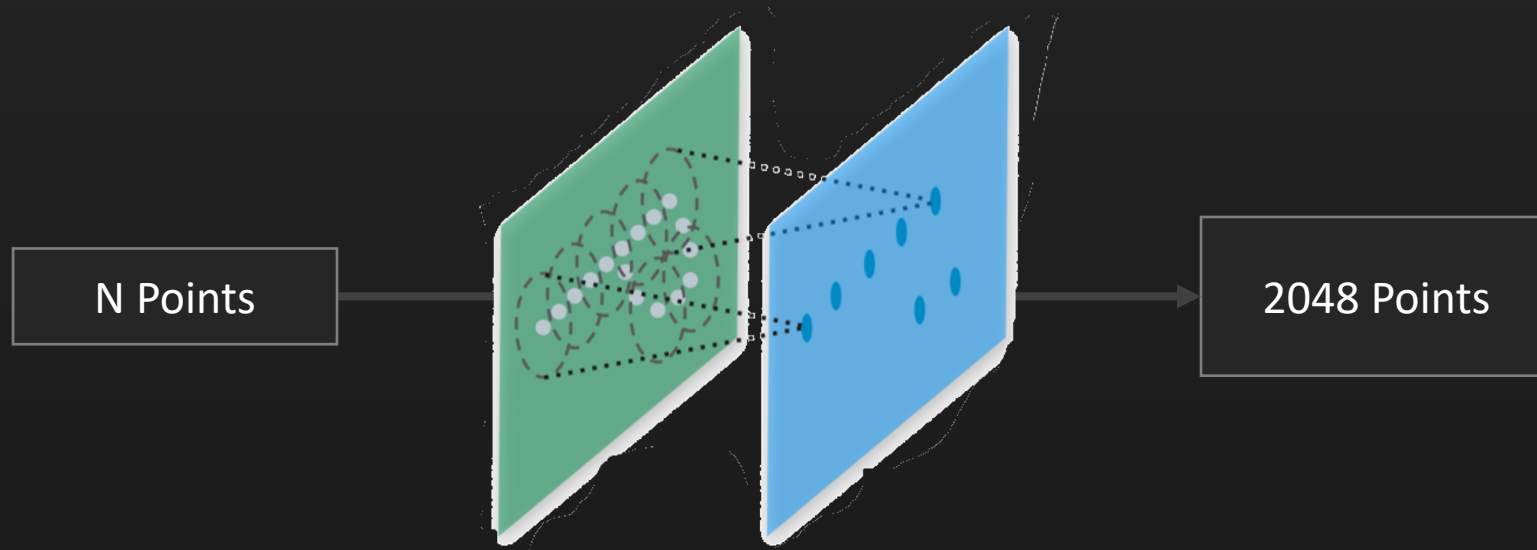
IV. Detection with Transformers



IV. Detection with Transformers



IV. Detection with Transformers



IV. Detection with Transformers

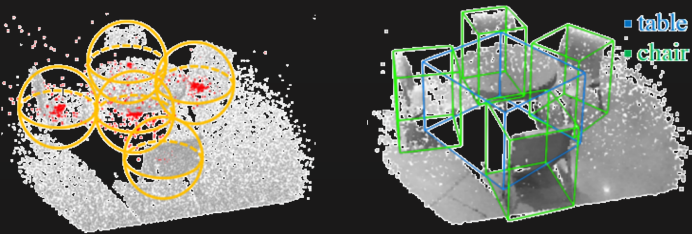
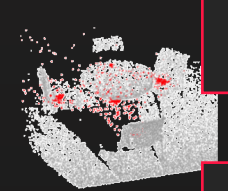
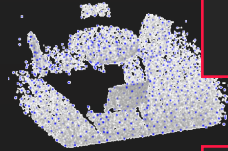
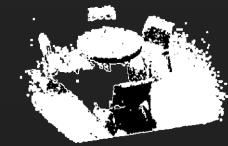
VoteNet

Point Cloud

PointNet++

Voting

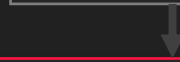
Proposal



3DETR

Point Cloud

Set Abstraction



IV. Detection with Transformers

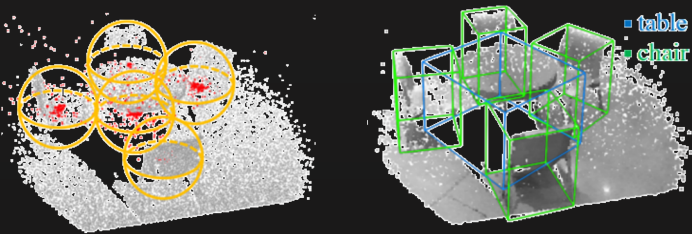
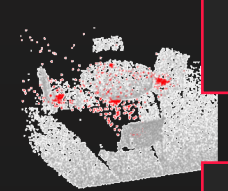
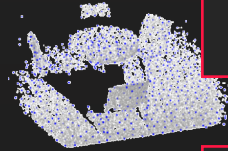
VoteNet

Point Cloud

PointNet++

Voting

Proposal



3DETR

Point Cloud

Set Abstraction



IV. Detection with Transformers

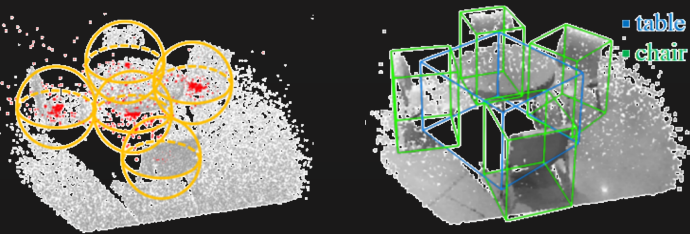
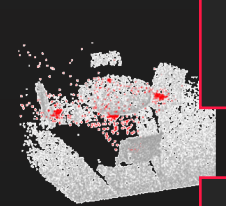
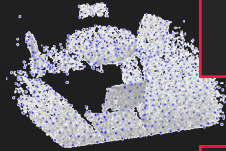
VoteNet

Point Cloud

PointNet++

Voting

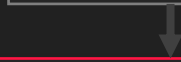
Proposal



3DETR

Point Cloud

Set Abstraction



Group-Free-3D

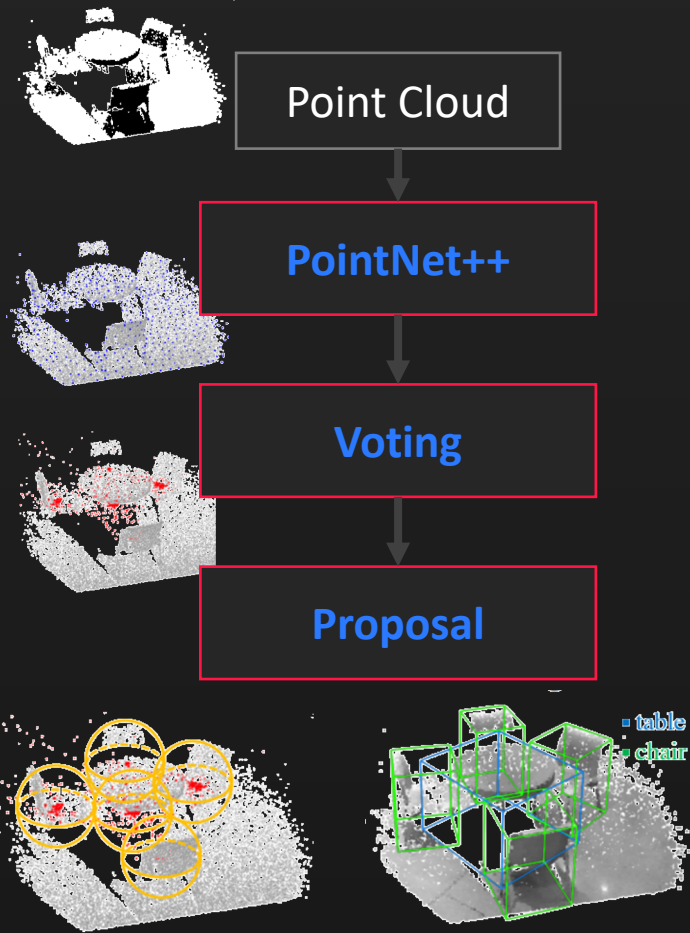
Point Cloud

Pointnet++

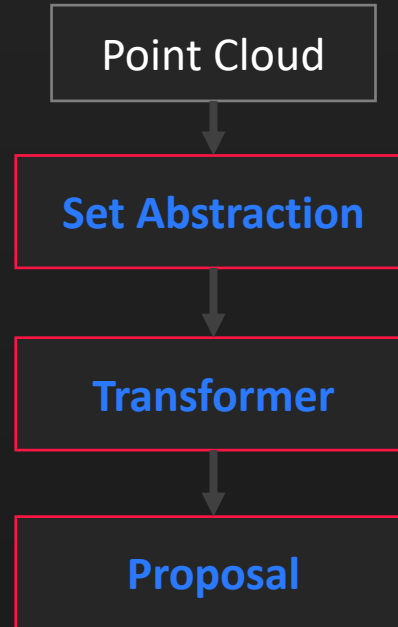


IV. Detection with Transformers

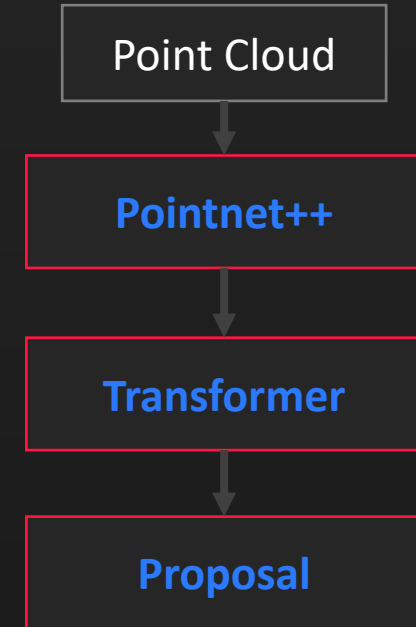
VoteNet



3DETR



Group-Free-3D



IV. Detection with Transformers

VoteNet

- Vote grouping is an issue!
- Cluster radius
- NMS reliance

3DETR

Group-Free-3D

IV. Detection with Transformers

VoteNet

- Vote grouping is an issue!
- Cluster radius
- NMS reliance

3DETR

+No NMS

Group-Free-3D

+No NMS

IV. Detection with Transformers

VoteNet

- Vote grouping is an issue!
- Cluster radius
- NMS reliance

3DETR

- +No NMS
- +Predict with every decoder output

Group-Free-3D

- +No NMS
- +Predict with every decoder output

IV. Detection with Transformers

VoteNet

- Vote grouping is an issue!
- Cluster radius
- NMS reliance

3DETR

- +No NMS
- +Predict with every decoder output

Group-Free-3D

- +No NMS
- +Predict with every decoder output
- +Use learnable pos. embeddings

IV. Detection with Transformers

VoteNet

- Vote grouping is an issue!
- Cluster radius
- NMS reliance

3DETR

- +No NMS
- +Predict with every decoder output

Group-Free-3D

- +No NMS
- +Predict with every decoder output
- +Use learnable pos. embeddings
- +More efficient point sampling strategy

IV. Detection with Transformers

VoteNet

- Vote grouping is an issue!
- Cluster radius
- NMS reliance

3DETR

- +No NMS
- +Predict with every decoder output
- +Simplest, flexible

Group-Free-3D

- +No NMS
- +Predict with every decoder output
- +Use learnable pos. embeddings
- +More efficient point sampling strategy

IV. Detection with Transformers

VoteNet

[mAP@0.5](#): 39.9

3DETR

[mAP@0.5](#): 47

Group-Free-3D

[mAP@0.5](#): 49

IV. Detection with Transformers

VoteNet

[mAP@0.5](#): 39.9

1M Parameters

3DETR

[mAP@0.5](#): 47

7.3M Parameters

Group-Free-3D

[mAP@0.5](#): 49

14.5M Parameters

IV. Detection with Transformers

VoteNet

mAP@0.5: 39.9

1M Parameters

3DETR

mAP@0.5: 47

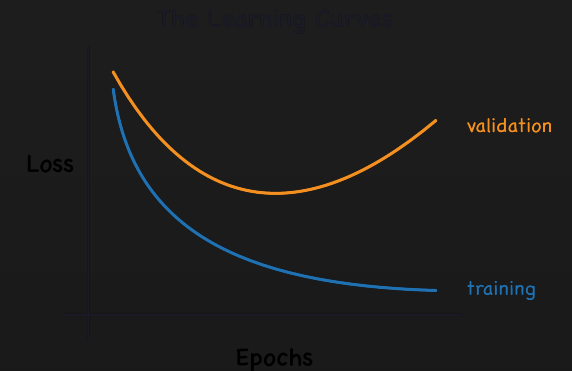
7.3M Parameters

Group-Free-3D

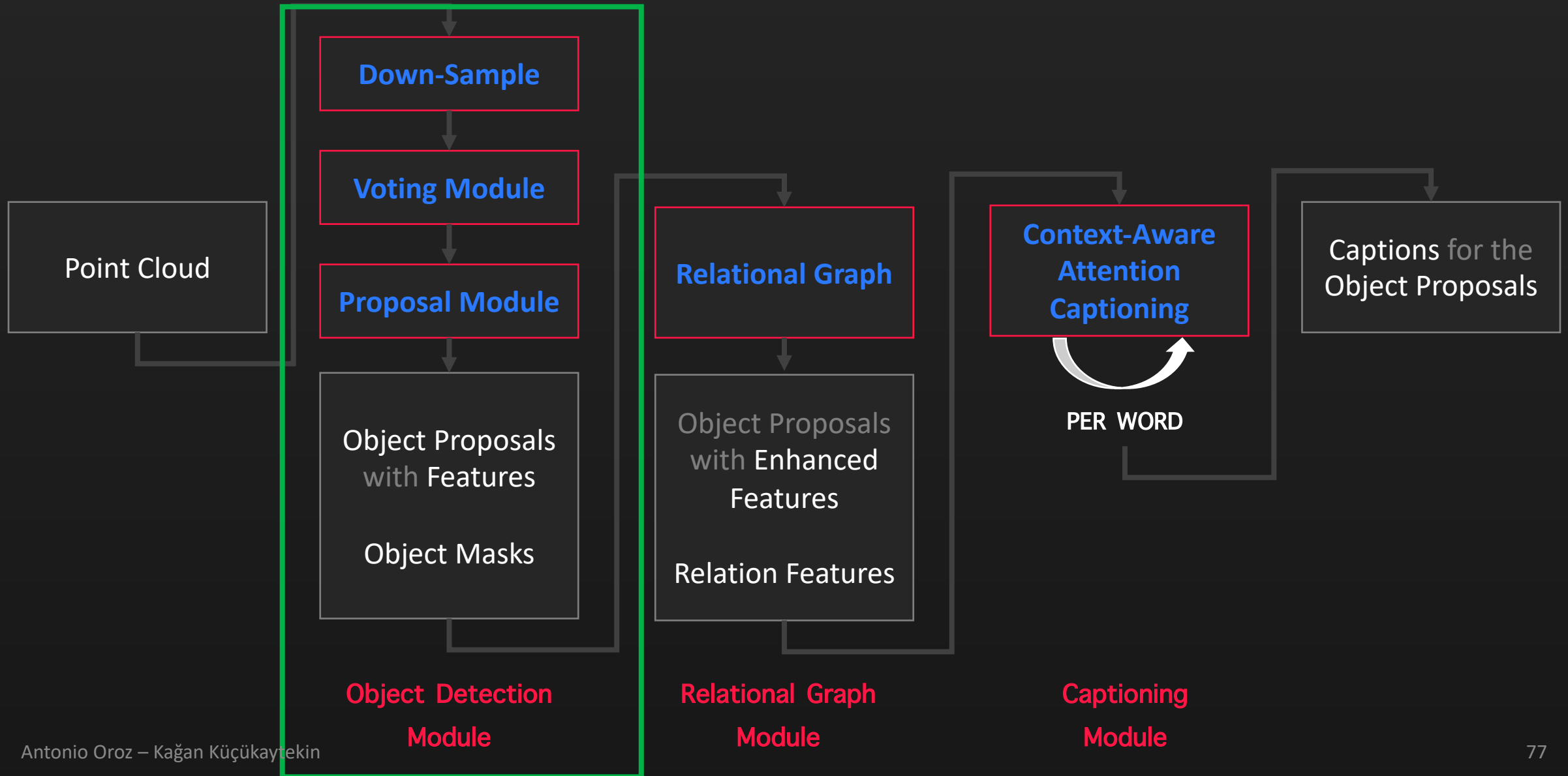
mAP@0.5: 49

14.5M Parameters

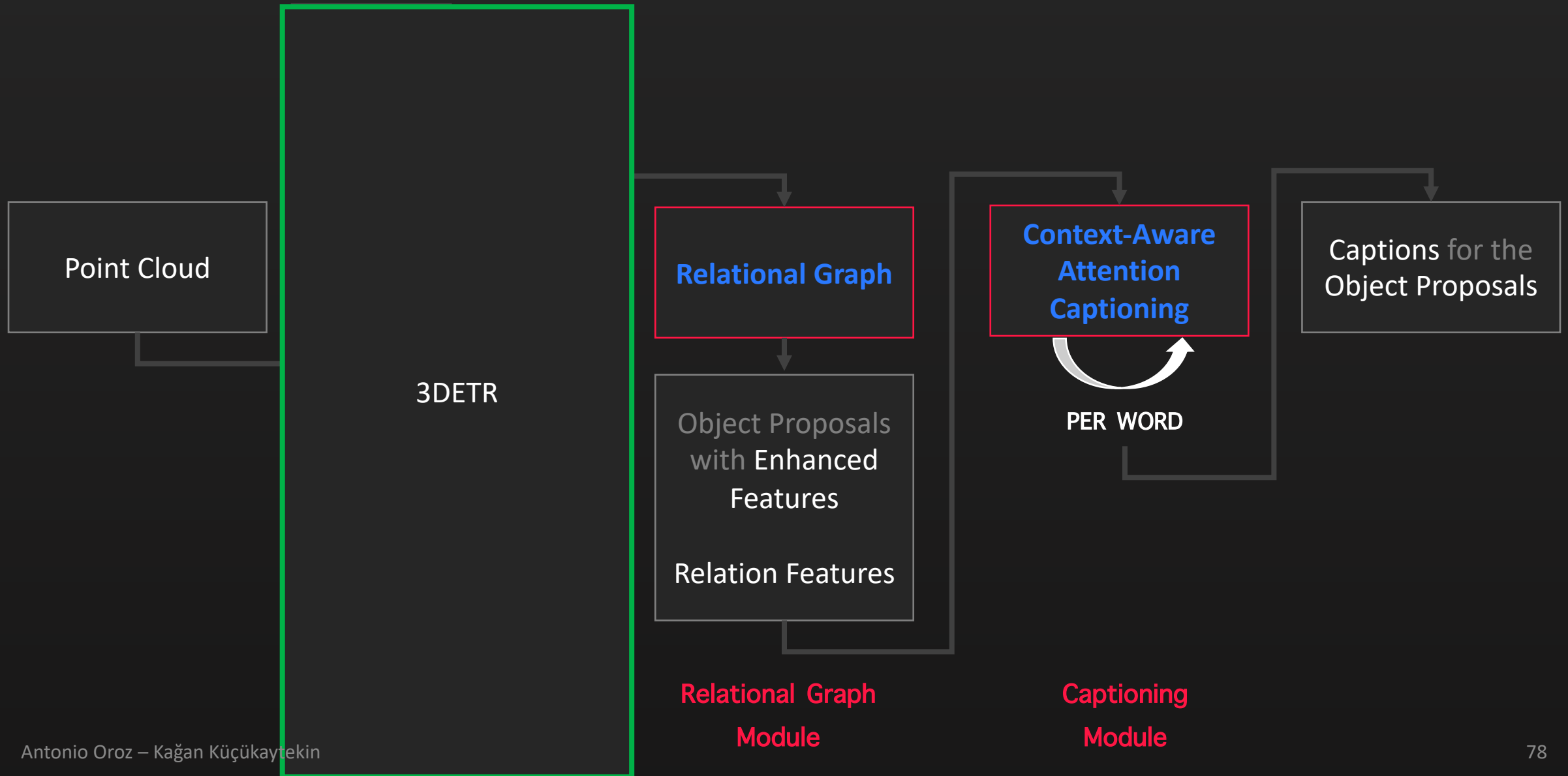
Because of data constraints, Group-Free-3D is more likely to overfit, so examine 3DETR



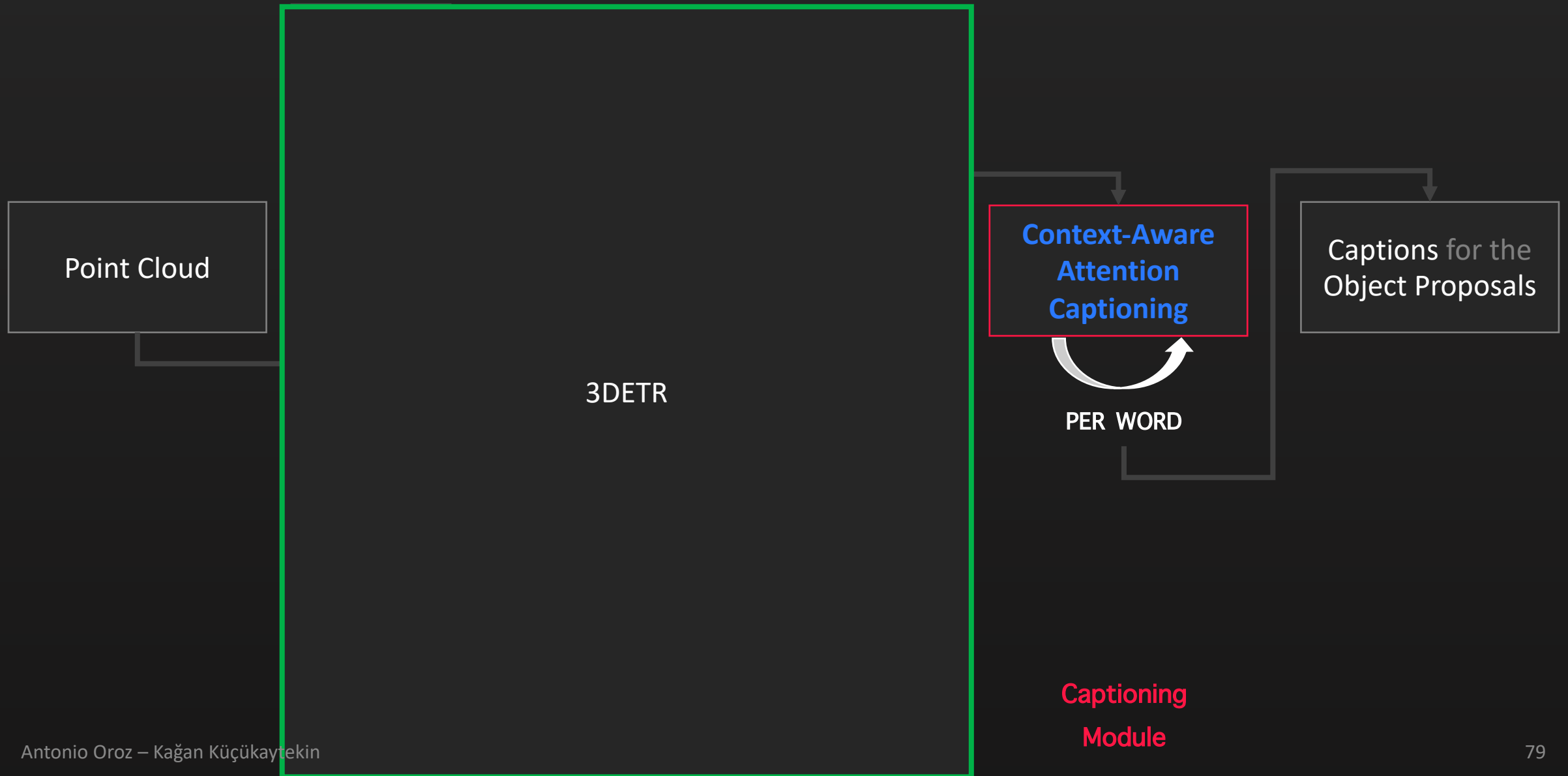
IV. Exploring Transformers for Detection Module



IV. Exploring Transformers for Detection Module



IV. Exploring Transformers for Detection Module



IV. Status

What we have done?

IV. Status

What we have done:

- Integrate 3DETR into the architecture

IV. Status

What we have done:

- Integrate 3DETR into the architecture
- Test end-to-end pipeline by overfitting to single sample without caption

IV. Status

What we have done:

- Integrate 3DETR into the architecture
- Test end-to-end pipeline by overfitting to single sample without caption

Possible next steps?

IV. Status

What we have done:

- Integrate 3DETR into the architecture
- Test end-to-end pipeline by overfitting to single sample without caption

Possible next steps:

- End-to-end overfit to small sample for whole task

IV. Status

What we have done:

- Integrate 3DETR into the architecture
- Test end-to-end pipeline by overfitting to single sample without caption

Possible next steps:

- End-to-end overfit to small sample for whole task
- Try transfer Learning with pre-trained 3DETR-m

IV. Status

What we have done:

- Integrate 3DETR into the architecture
- Test end-to-end pipeline by overfitting to single sample without caption

Possible next steps:

- End-to-end overfit to small sample for whole task
- Try transfer Learning with pre-trained 3DETR-m
- No promise! Ablation studies on our model is our Prio 1.

Scan2CapMMT

- I. Scan2CapMMT Recap
- II. Improving Scan2CapMMT
- III. Quantitative & Qualitative Results
- IV. Detection with Transformers
- V. Timeline until the Final Presentation

Scan2CapMMT

- I. Scan2CapMMT Recap
- II. Improving Scan2CapMMT
- III. Quantitative & Qualitative Results
- IV. Detection with Transformers
- V. Timeline until the Final Presentation

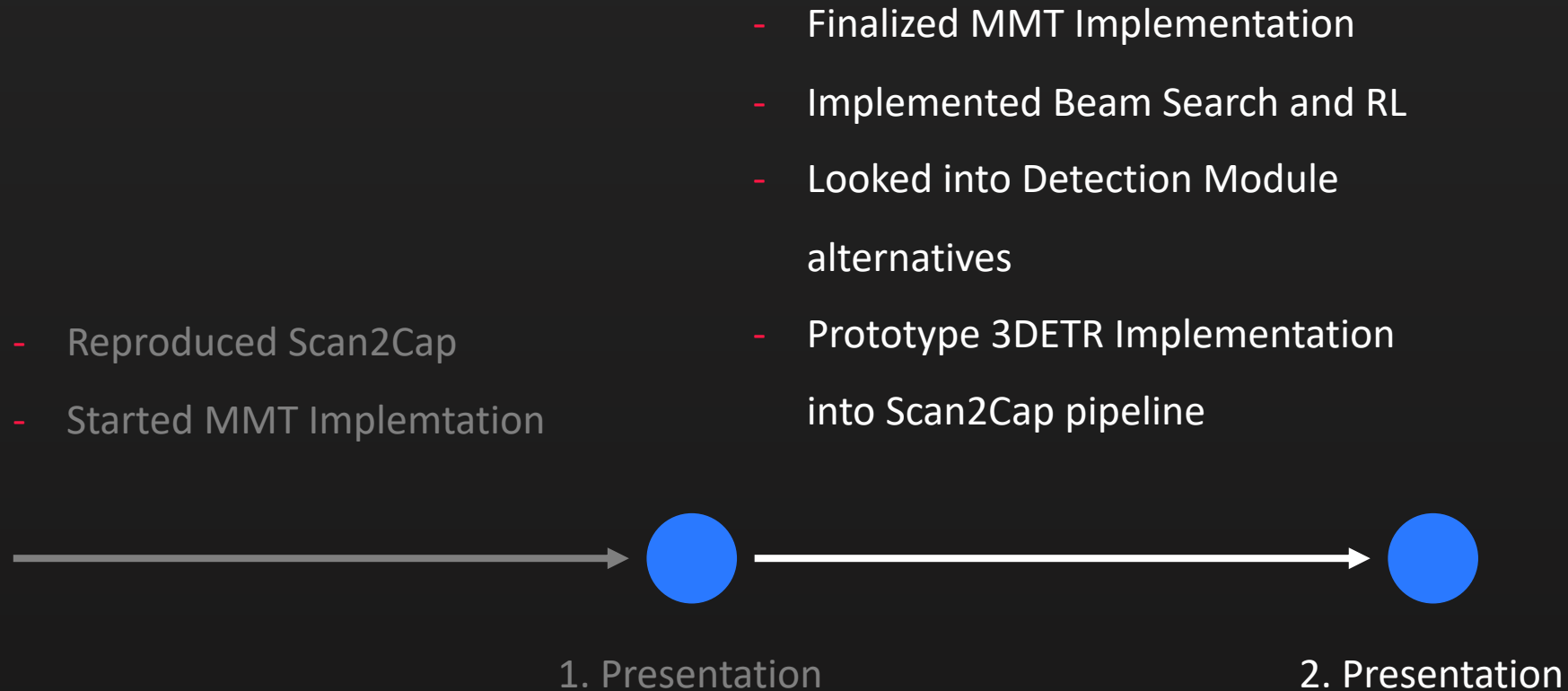
V. Timeline until Final Presentation

- Reproduced Scan2Cap
- Started MMT Implementation

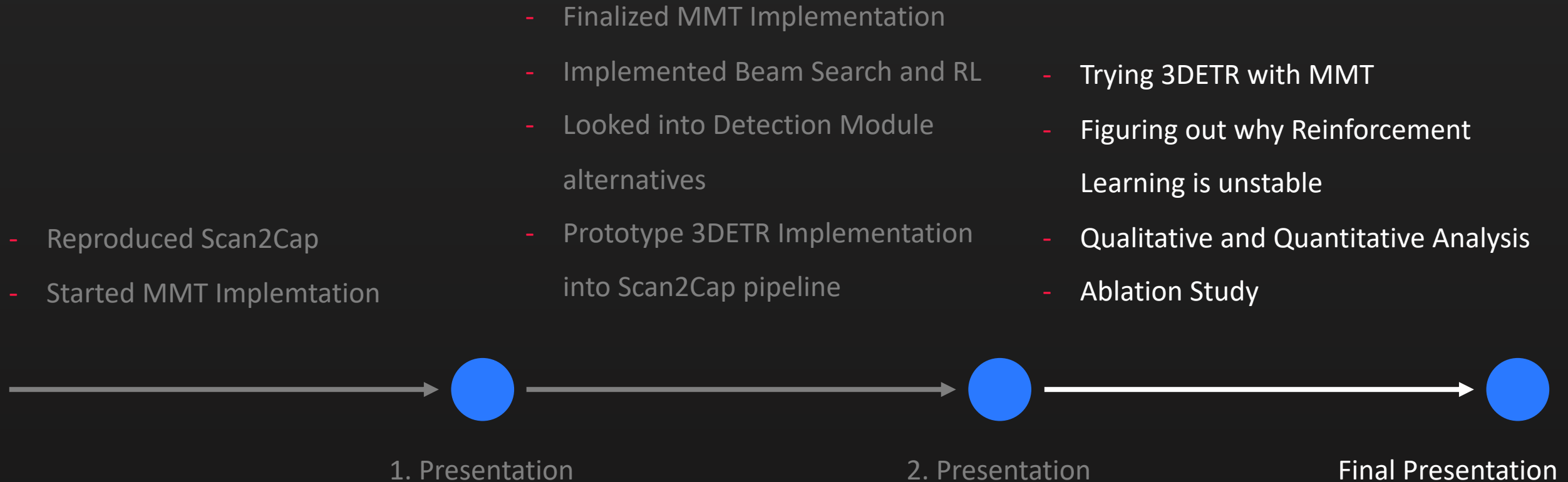


1. Presentation

V. Timeline until Final Presentation



V. Timeline until Final Presentation



Scan2CapMMT

- I. Scan2CapMMT Recap
- II. Improving Scan2CapMMT
- III. Quantitative & Qualitative Results
- IV. Detection with Transformers
- V. Timeline until the Final Presentation

Scan2CapMMT

- I. Scan2CapMMT Recap
- II. Improving Scan2CapMMT
- III. Quantitative & Qualitative Results
- IV. Detection with Transformers
- V. Timeline until the Final Presentation

**THANK YOU FOR
YOUR ATTENTION :D**