# MMC Transformer: Multiscale Multigrid Comparator Transformer for Few-Shot Video Segmentation

Mennatullah Siam, Konstantinos G. Derpanis, Richard Wildes

**Abstract:** ...

## Paper Decision

**Decision:** Accept

## Meta Review of Paper7524 by Area Chair Daity

## Author Rebuttal Acknowledgement by Paper7524 Reviewer DGqR

## Author Rebuttal Acknowledgement by Paper7524 Reviewer HuqG

## Author Rebuttal Acknowledgement by Paper7524 Reviewer ZBpu

## No Reviewer Responses

## Strengths + Novelty

## Official Review of Paper7524 by Reviewer HuqG

**Summary:**

**Strengths:**

**Weaknesses:**

### R1 response

## Official Review of Paper7524 by Reviewer DGqR

**Summary:**

**Strengths And Weaknesses:**

**Limitations:**

### R2 response

## Official Review of Paper7524 by Reviewer 97JM

**Summary:**

**Strengths And Weaknesses:**

**Limitations:**

### R3 response

### Response

## Official Review of Paper7524 by Reviewer ZBpu

**Summary:**

**Strengths And Weaknesses:**

**Limitations:**

### R4 response

### Response