

---

# Diffusion Model for Super Resolution of Ship Images

Student

**Andrea Mazzitelli - 1835022**

Advisor

**Prof. Danilo Comminiello**

Engineering in Computer Science

Dipartimento di Ingegneria informatica, automatica e gestionale

Facoltà di Ingegneria dell'informazione, informatica e statistica

Anno Accademico 2022-2023

Sapienza University of Rome

Co-Advisor

**Ing. Alessandro Nicolosi**



---

In Collaboration with



---

# Table of contents

1

Introduction

3

Chosen Approach

5

Training and  
Results

2

Baseline

4

Dataset

6

Conclusions and  
Future Works

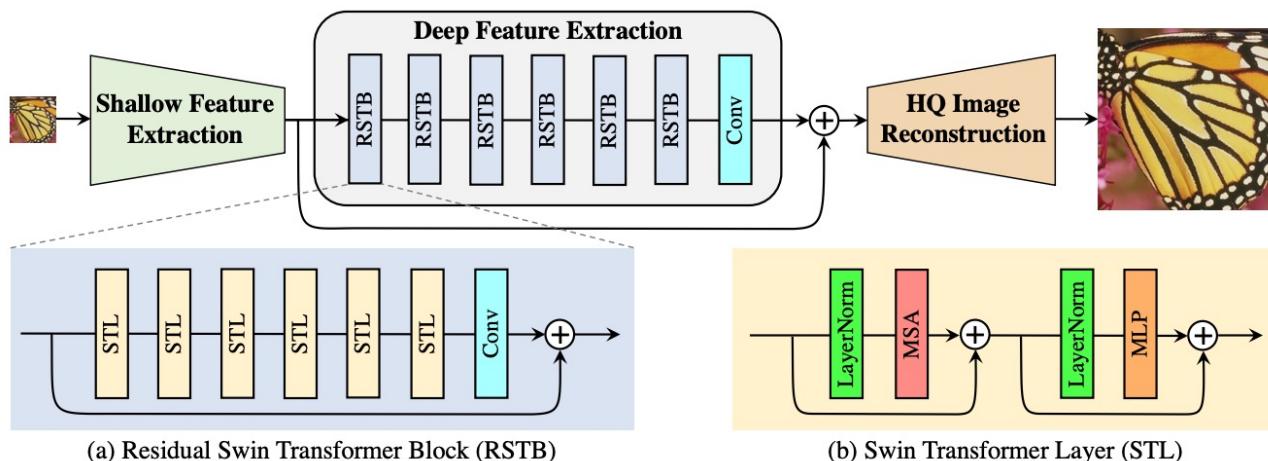




# Background - Image Super Resolution



# Baseline - SwinIR

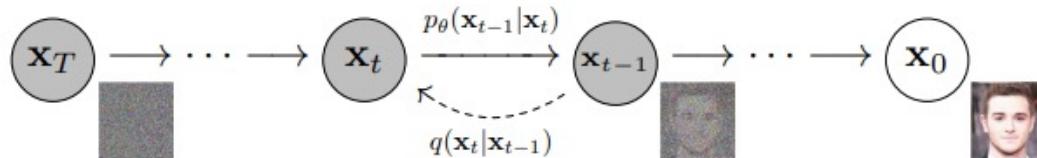




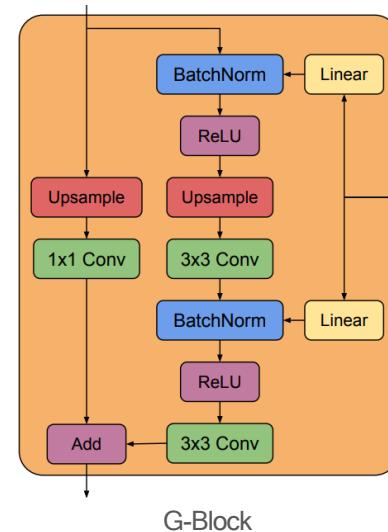
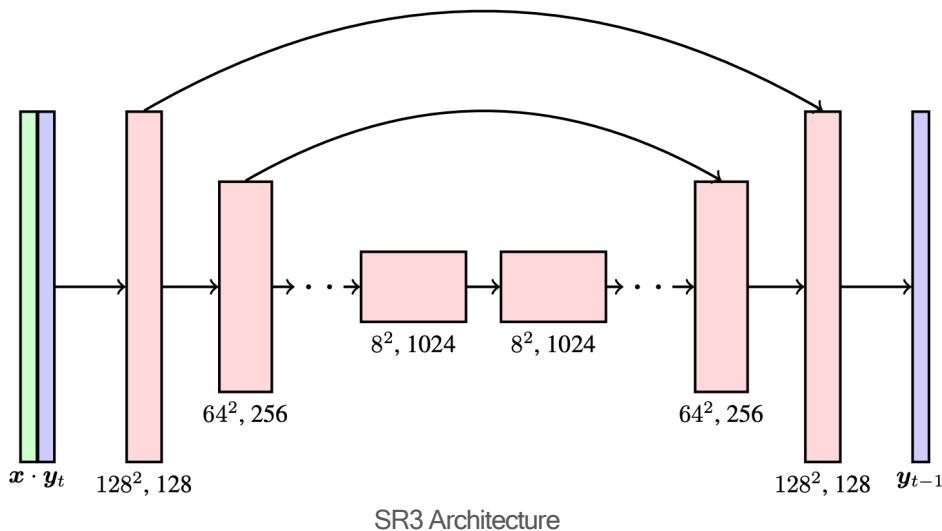
# Chosen Approach - Diffusion Model

Diffusion Models leverage a diffusion process to iteratively transform an initial distribution into a target one. This is achieved using a **Forward Process** and a **Reverse Process**.

By including a **Conditioning Element** in the process is possible to guide the diffusion process in generating an output that have specific characteristics.

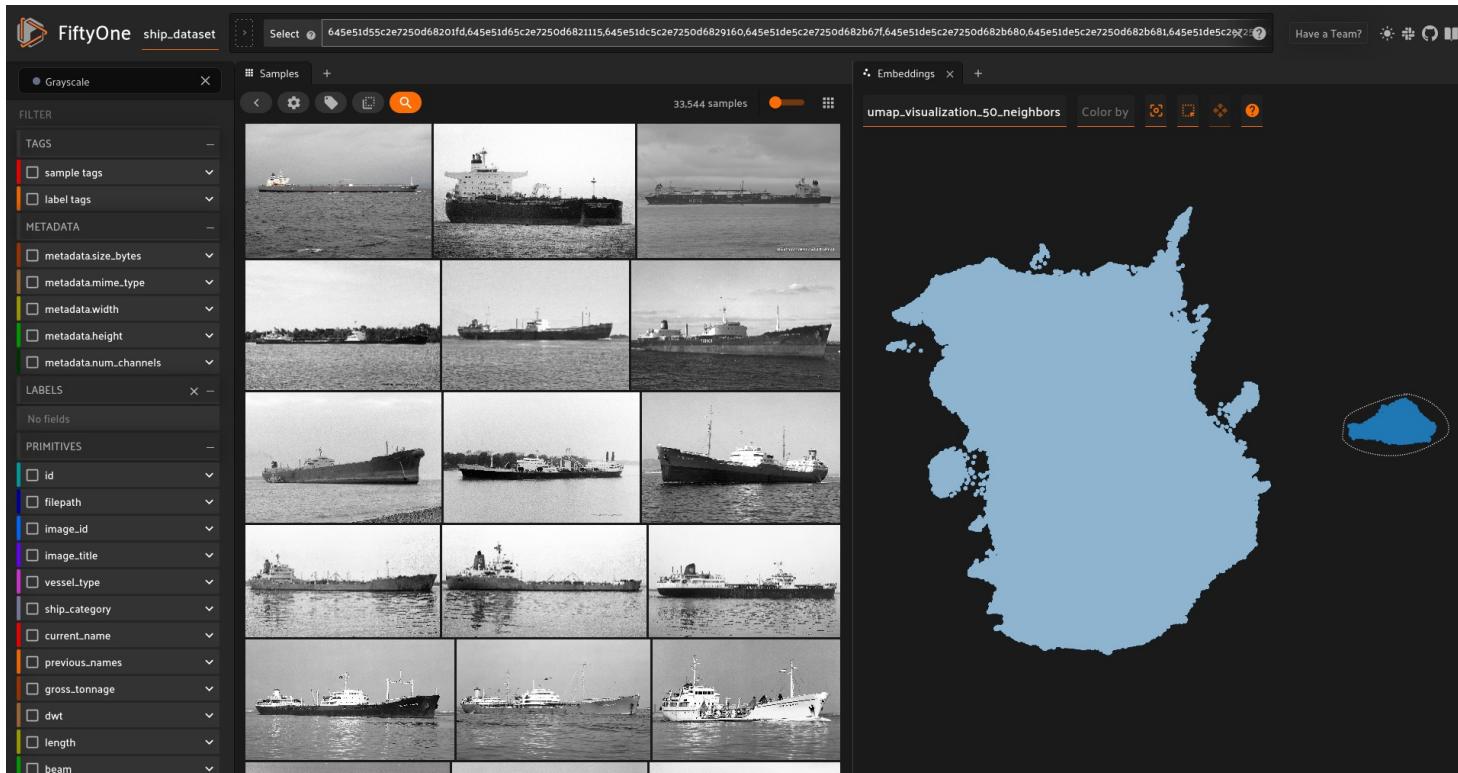


# Super Resolution via Repeated Refinement - SR3



# Dataset





---

# Training Details

**TASK :** 8x Super Resolution

64x64 → 512x512

**Devices :** 2 NVIDIA RTX A6000, 48GB



- **SwinIR:**

- 12 million parameters
- 500k iterations
- batch size :32
- 10 days
- inference time: few seconds

- **SR3:**

- 155 million parameters
- 1M iterations
- batch size: 16
- 17 days
- inference time: 2:14



## Results - Objective Metrics

	PSNR (db)	SSIM	FID
<b>SwinIR</b>	<b>23.80</b>	<b>0.6805</b>	50.79
<b>SR3</b>	21,07	0.57093	<b>16.59</b>

PSNR - Peak Signal to Noise Ratio, SSIM - Structural Similarity Index Measures, FID - Fréchet Inception Distance

SwinIR outperforms SR3 in both PSNR and SSIM but not in FID. This may appear to be somewhat counterintuitive but there is a reason for it.





# Results - Subjective Metric

2AFC - 2 Alternative Force Choice



INPUT



SwinIR



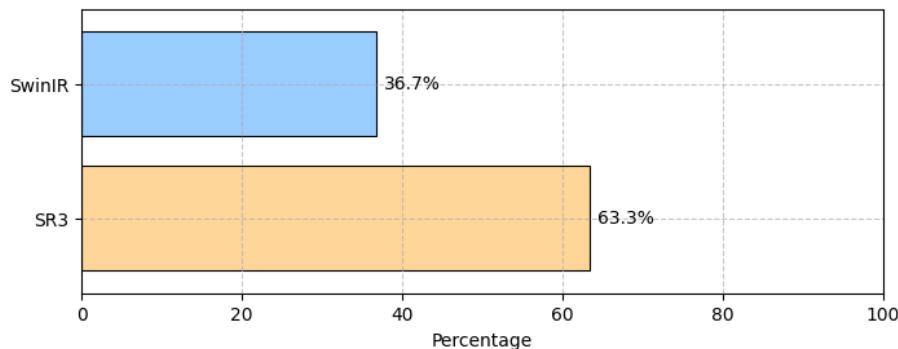
SR3



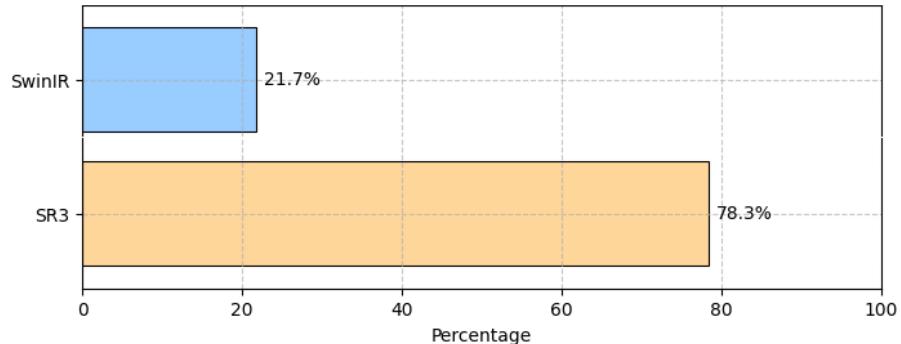


## Results - Subjective Metrics contd.

Which is the best high-resolution version of the given low-resolution image?



Which is the most realistic version of the given low-resolution image?



---

# Results - Visual Analysis



INPUT



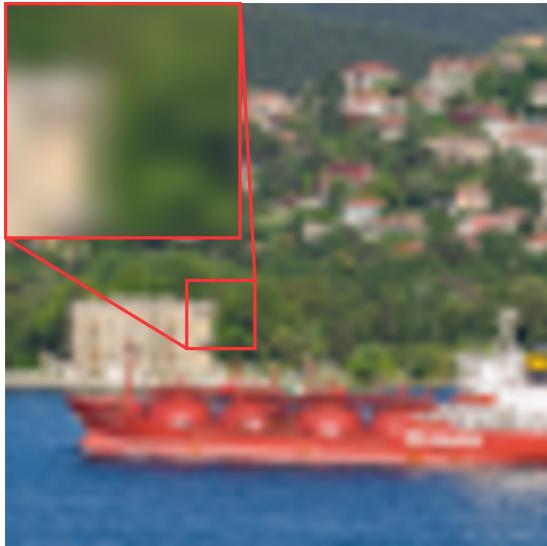
SwinIR



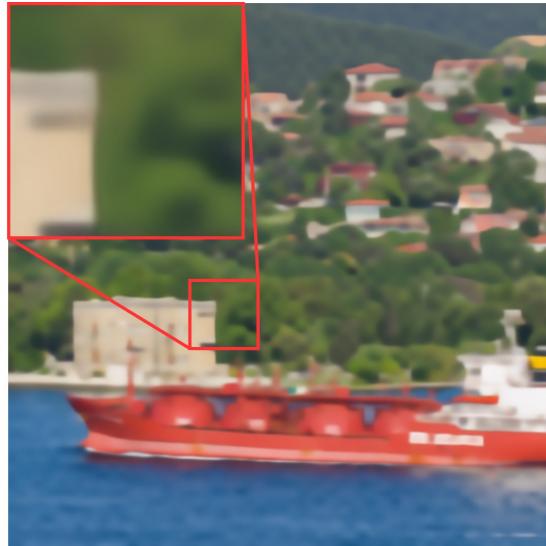
SR3



# Results - Visual Analysis



INPUT



SwinIR



SR3



---

# Results - Visual Analysis



INPUT



SwinIR



SR3

---

# Results - Visual Analysis



INPUT



SwinIR



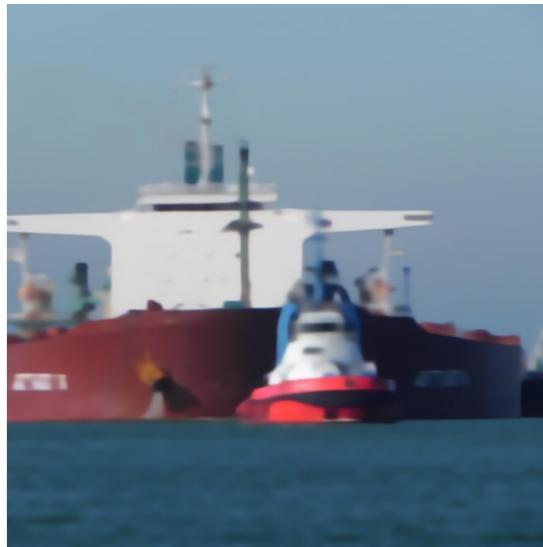
SR3

---

# Results - Visual Analysis



INPUT



SwinIR



SR3

---

# Results - Visual Analysis



INPUT



SwinIR



SR3

---

# Results - Visual Analysis



INPUT



SwinIR



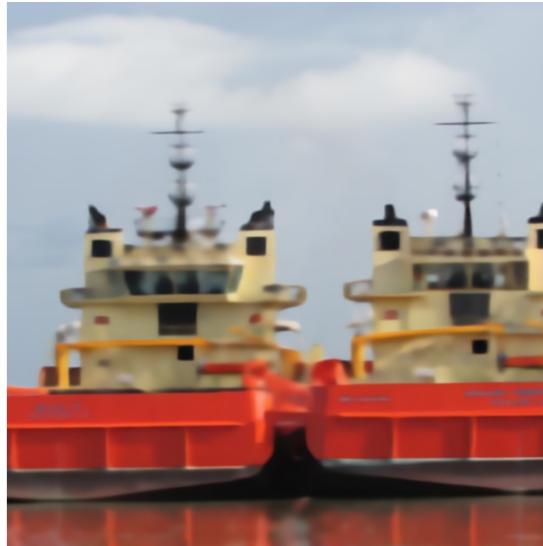
SR3

---

# Results - Visual Analysis



INPUT



SwinIR



SR3



## Conclusions and Future Works

- SwinIR : generates results that are perceived as worse by the human eye but inference is fast
- SR3 : produces higher quality images and more realistic to the human eye but inference is long

**TRADE OFF : Image Quality vs Inference Speed**

**Future Works : reduce inference time of the SR3 model**





# Thank you for the attention!

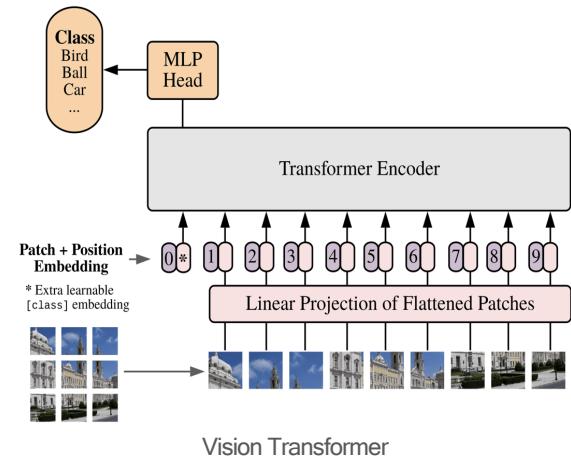


# Transformers for Computer Vision

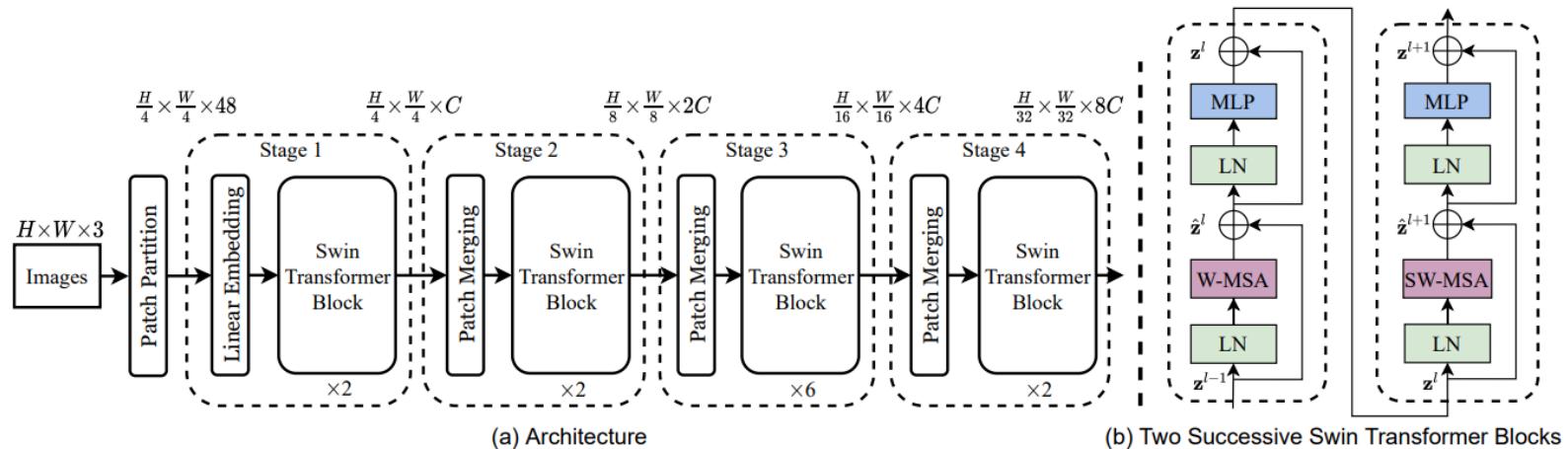
Transformers are encoder-decoder architecture built to process sequence of tokens. So by default they are impractical for computer vision task due to the high number of tokens needed to represent an image.

The **Vision Transformer** solved this issue by considering patches of the image as tokens, reducing the length of the sequence representing the image. Still, ViT is limited to process the image on a single scale

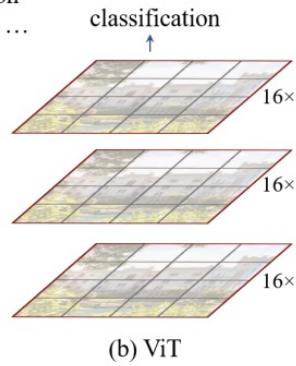
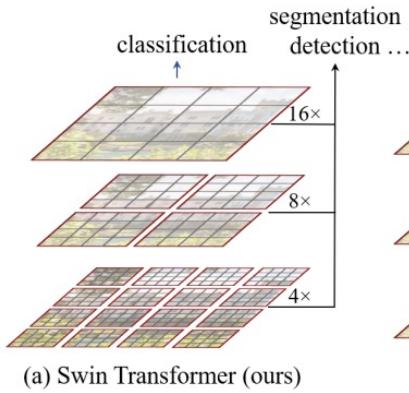
So, the **Swin Transformer** introduce a hierarchical shift-based window mechanism with local self-attention to process images at multiple scales.



# Swin Transformer



# Swin Transformer



: *shifted window approach* :

A local window to  
perform self-attention  
A patch





# SR3 Generation Process

