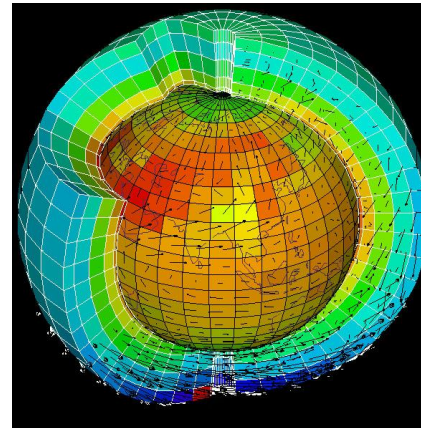


GENERATIVE ML IN ATMOSPHERIC MODELS

Varshini Thavakumar

PROBLEM BACKGROUND

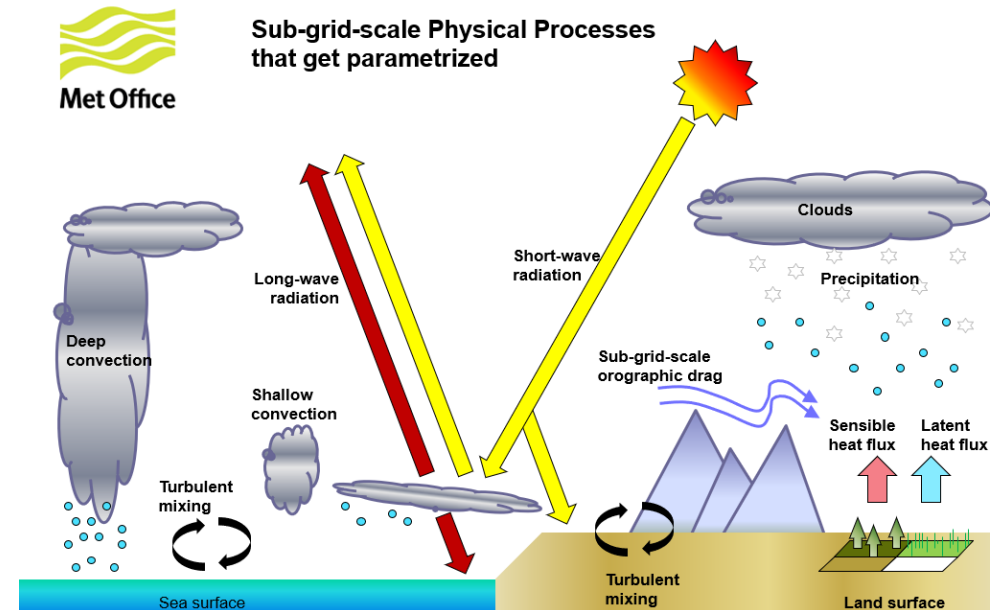


$$\frac{\partial \theta}{\partial t} = l - \frac{u}{a \cos \phi} \frac{\partial \theta}{\partial \lambda} - \frac{v}{a} \frac{\partial \theta}{\partial \phi} - w \frac{\partial \theta}{\partial z}$$

$$\frac{\theta(t+1) - \theta(t)}{\Delta t} = l - \frac{u}{a \cos \phi} \frac{\theta(i+1) - \theta(i-1)}{\Delta \lambda} - \frac{v}{a} \frac{\theta(j+1) - \theta(j-1)}{\Delta \phi} - w \frac{\theta(k+1) - \theta(k-1)}{\Delta z}$$

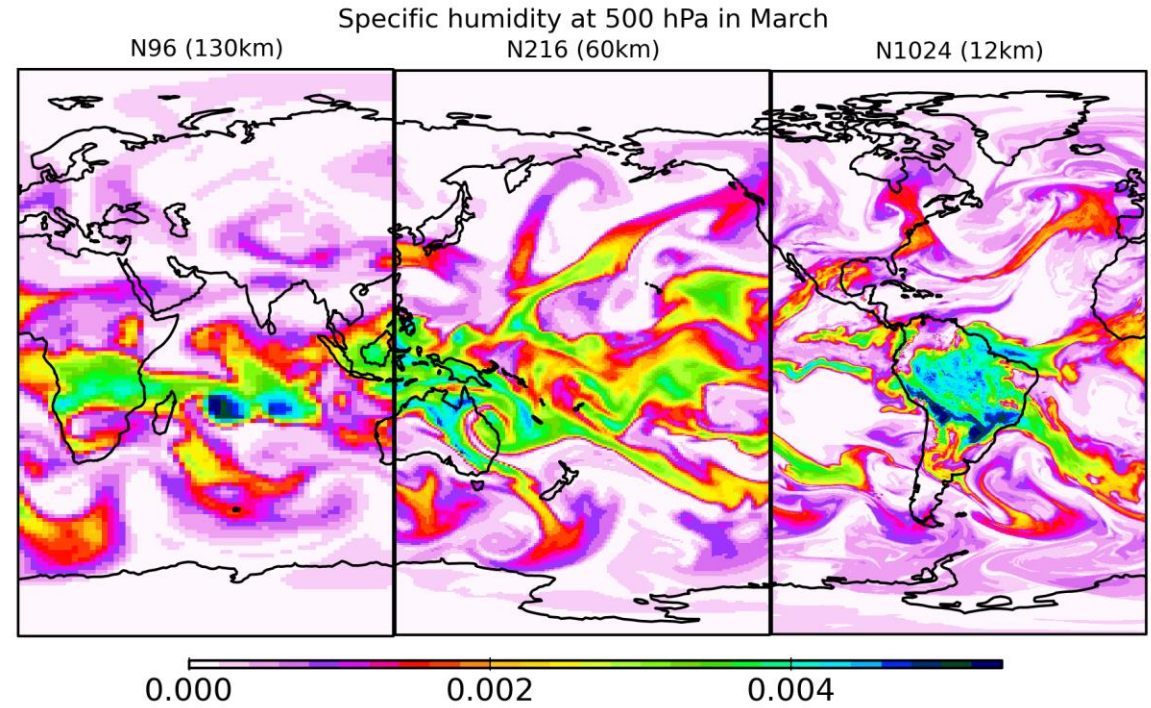
$$\theta(t+1) = \theta(t) + \Delta t \left\{ l - \frac{u}{a \cos \phi} \frac{\theta(i+1) - \theta(i-1)}{\Delta \lambda} - \frac{v}{a} \frac{\theta(j+1) - \theta(j-1)}{\Delta \phi} - w \frac{\theta(k+1) - \theta(k-1)}{\Delta z} \right\}$$

3-dimensional arrays representing wind, temperature, pressure, water vapour, clouds etc at all latitudes and longitude at a range of heights



Generally, these processes happen on scales much smaller than the grid, and so the Met Office parametrise them.

PROBLEM BACKGROUND

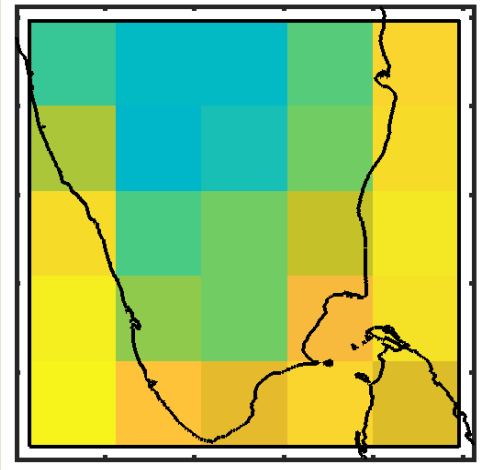


Using high resolution grids is limited by computational power

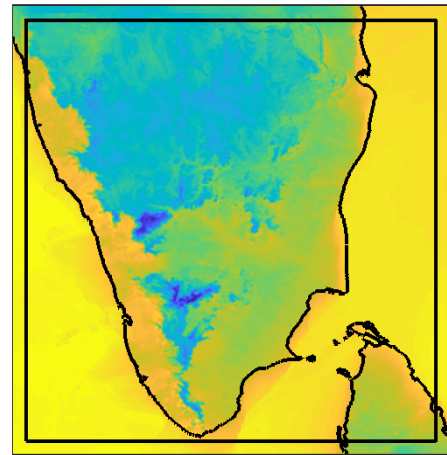
Halving grid length ~ **16 times more expensive !**

PROBLEM OBJECTIVE

(1) Given a low resolution grid, can we generate samples matching the expected sub-grid distribution?



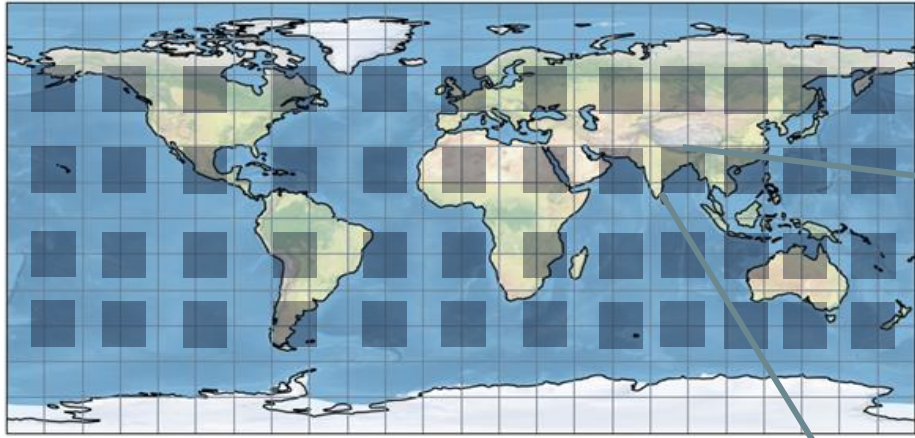
(2) And if we can, what happens when this is fed into the parametrisation schemes?



How?

Use Generative ML

DATA

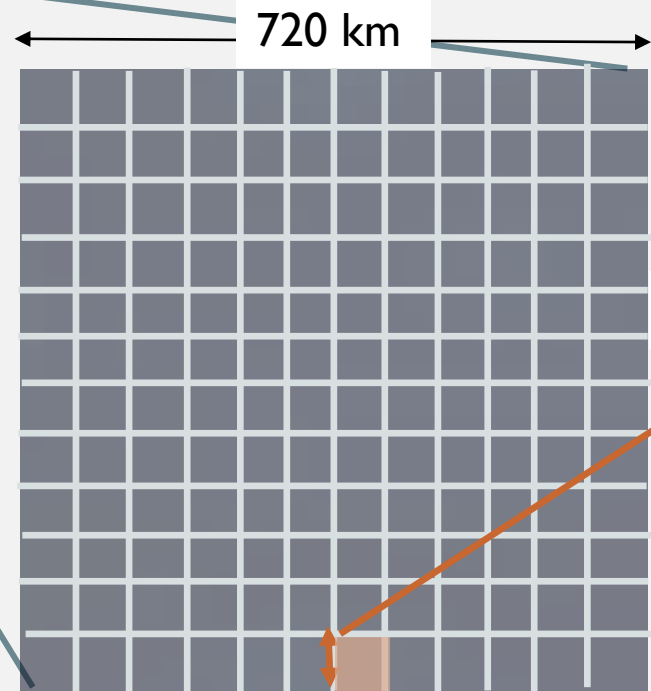


Focus on **one simulation only** with data in one 7° by 7° region with the following:

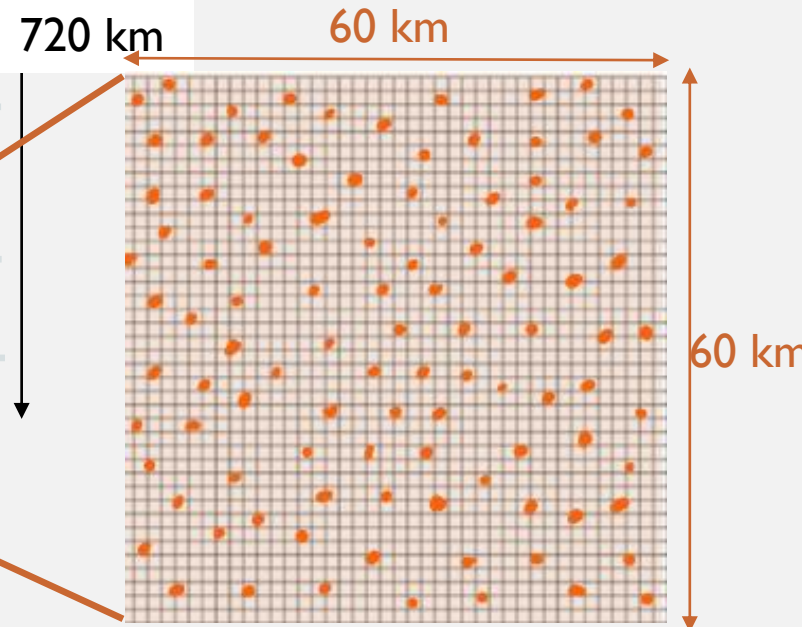
- Land-sea fraction
- Mean Orography
- Std dev Orography
- Temperature
- Pressure
- Specific Humidity

Data available hourly for whole month of January 2020.

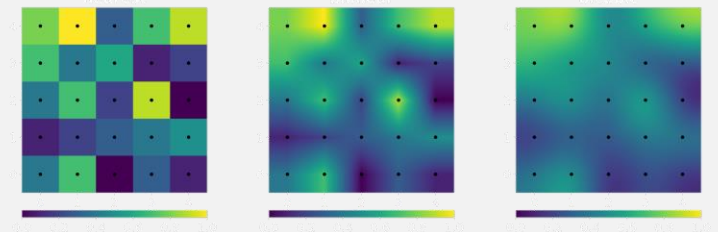
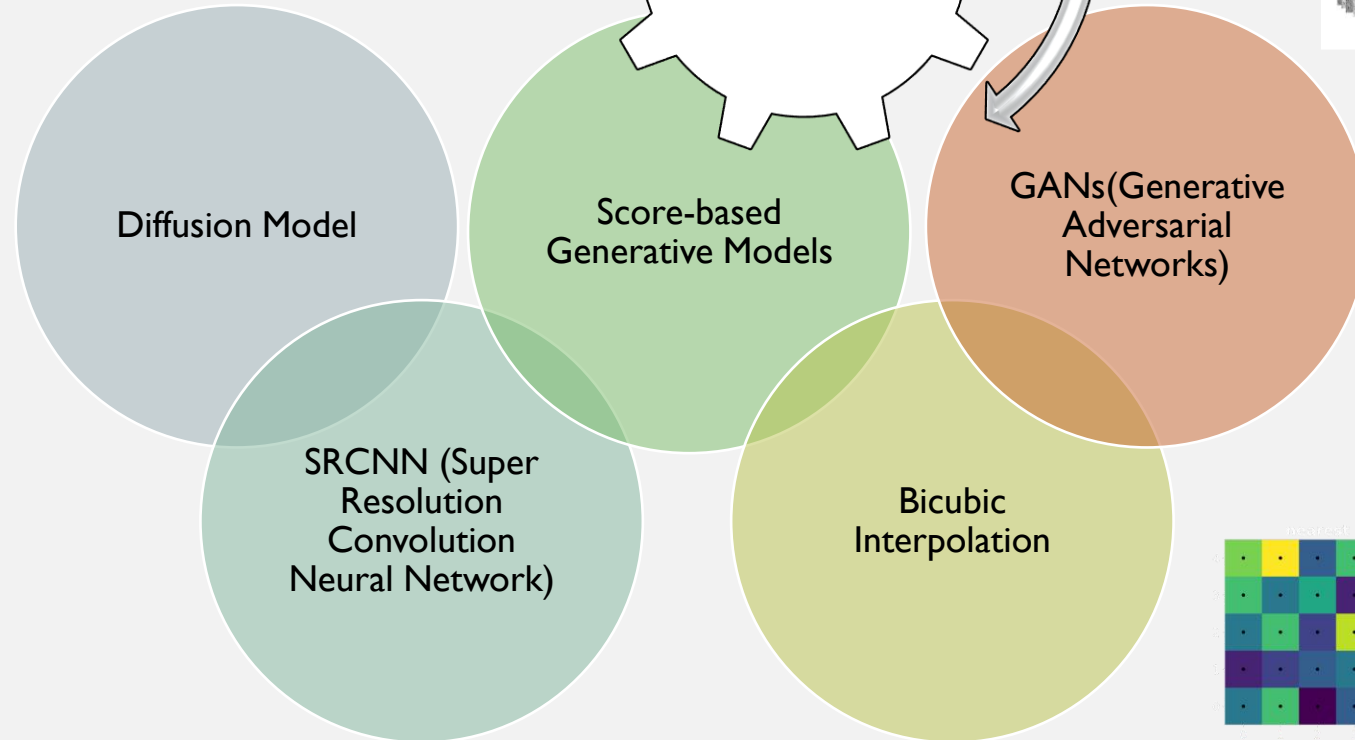
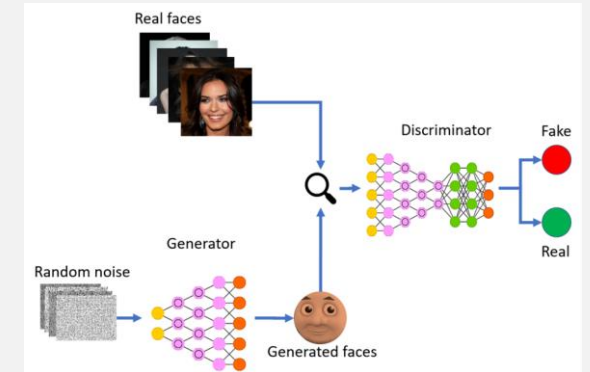
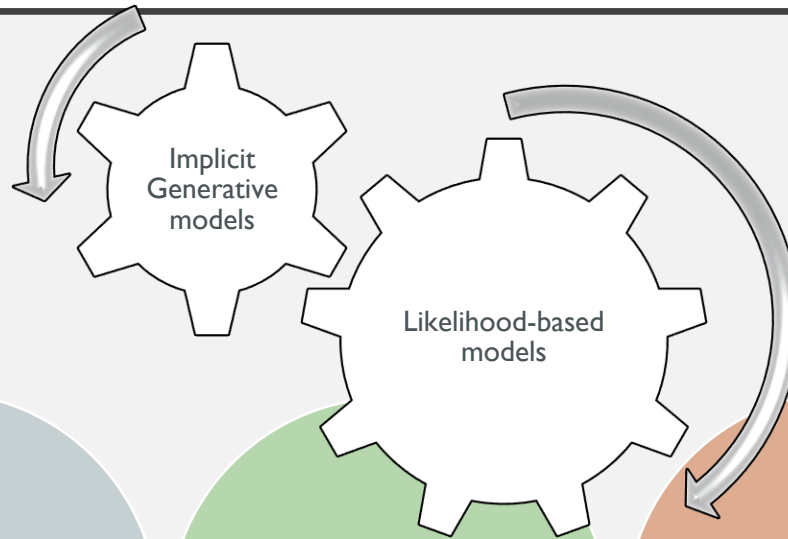
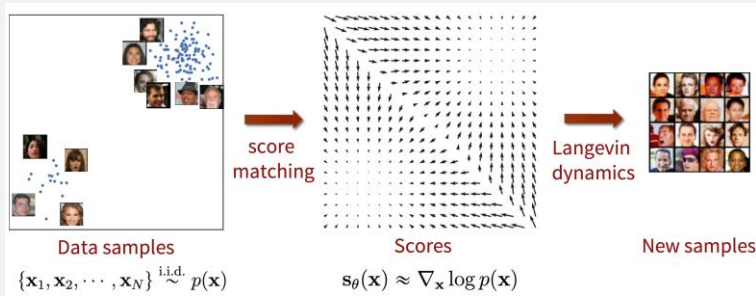
Met Office have ~5TB of data: 80 high resolution simulations of atmosphere (1.5km grid length) scattered all around the globe. (2% of the world)



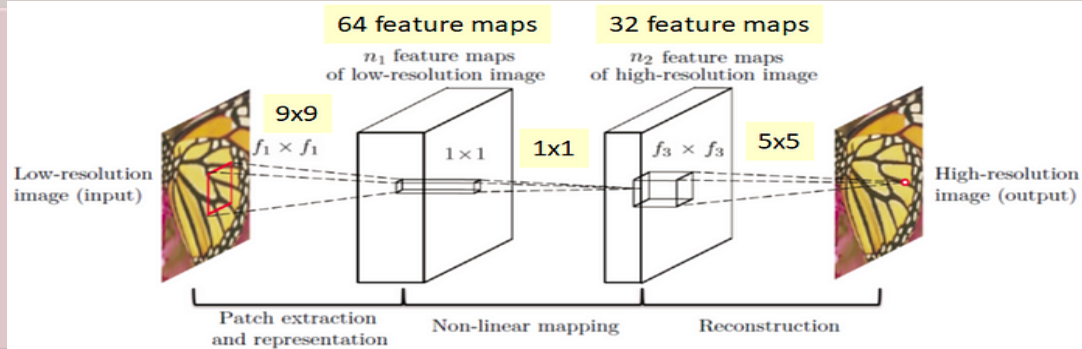
Within each 60km grid box, 100 random samples of 1.5km boxes given.



GENERATIVE MODELLING TECHNIQUES



SRCNN



Loss function

Average of mean square error for n training samples.

$$L(\Theta) = \frac{1}{n} \sum_{i=1}^n \|F(\mathbf{Y}_i; \Theta) - \mathbf{X}_i\|^2$$

1) Patch Extraction and Representation.

Low-resolution input is upscaled to the desired size using **bicubic interpolation**. First layer performs a standard convolution with ReLU to get

$$F_1(\mathbf{Y}) = \max(0, W_1 * \mathbf{Y} + B_1)$$

2) Non-linear Mapping.

Uses a 1x1 convolution. Maps low-resolution vector to high resolution vector

$$F_2(\mathbf{Y}) = \max(0, W_2 * F_1(\mathbf{Y}) + B_2)$$

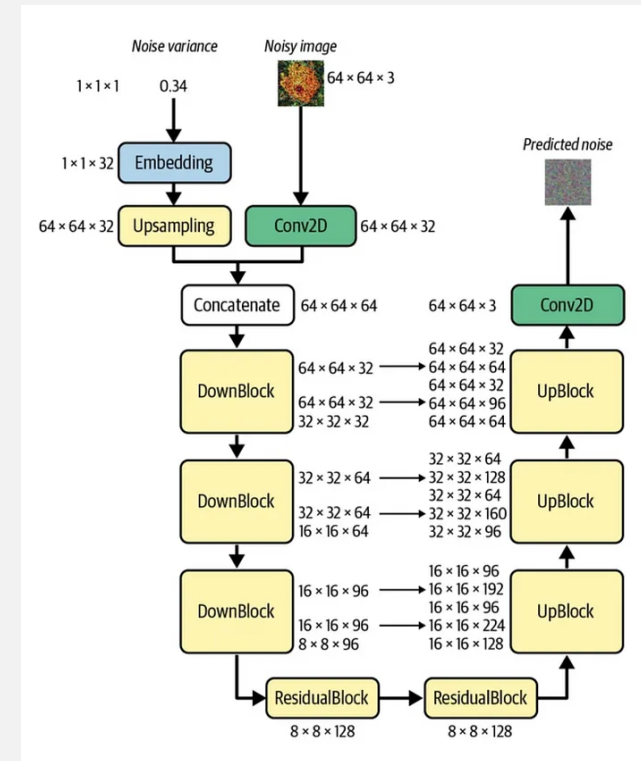
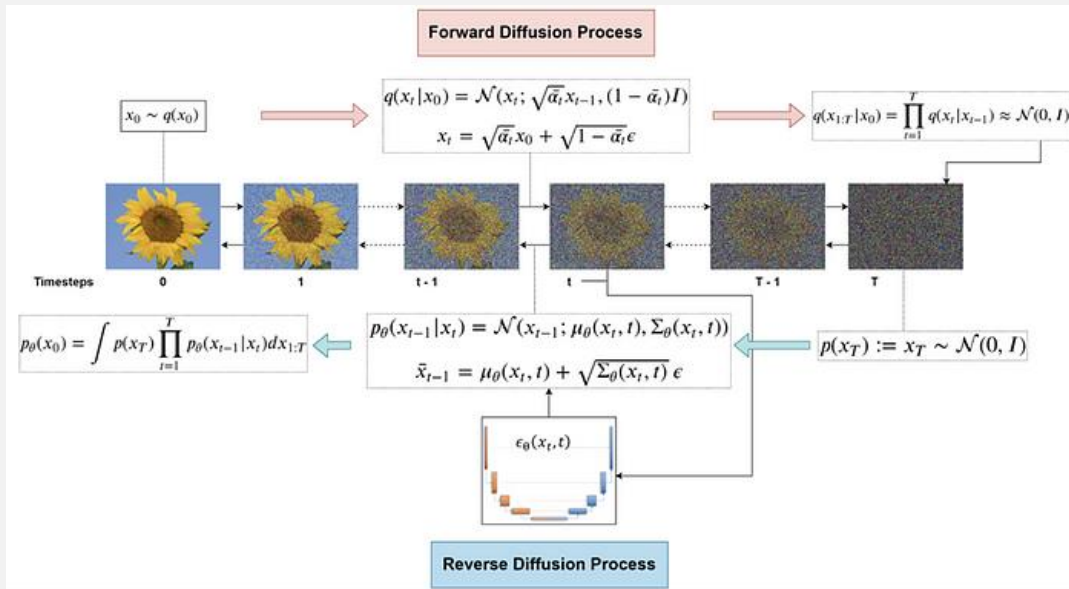
3) Reconstruction.

Convolution again to reconstruct the image.

$$F(\mathbf{Y}) = W_3 * F_2(\mathbf{Y}) + B_3$$

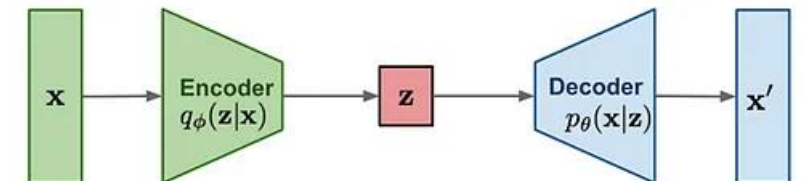
DIFFUSION MODELS

Diffusion models reconstruct images from noise. They slowly apply Gaussian noise to the training images, then a neural network is trained to predict the noise at each time step. Each time, noise is predicted, this is subtracted from the image until a realistic image emerges.



For resolution tasks, diffusion models would include a U-net architecture and possibly a Variational Autoencoder (VAE).

VAE: maximize ELBO.



TIMELINE, EXPECTED OUTCOME & COMPUTATIONAL REQUIREMENTS

Review of previous works, methods and techniques used. Finding code on github/Kaggle/papers. (2 weeks)



General literature review. (1.5 weeks)



In depth review of predictive techniques. (1.5 weeks)



In depth review of applicable machine learning techniques. (1.5 weeks)



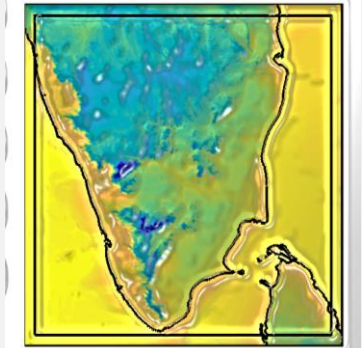
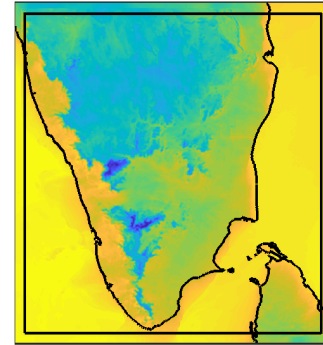
Model development (3 weeks)



Testing and performance measurement. (2 weeks)



Optimisation. (3 weeks)



Using our generated high resolution samples, we compare to the originals.

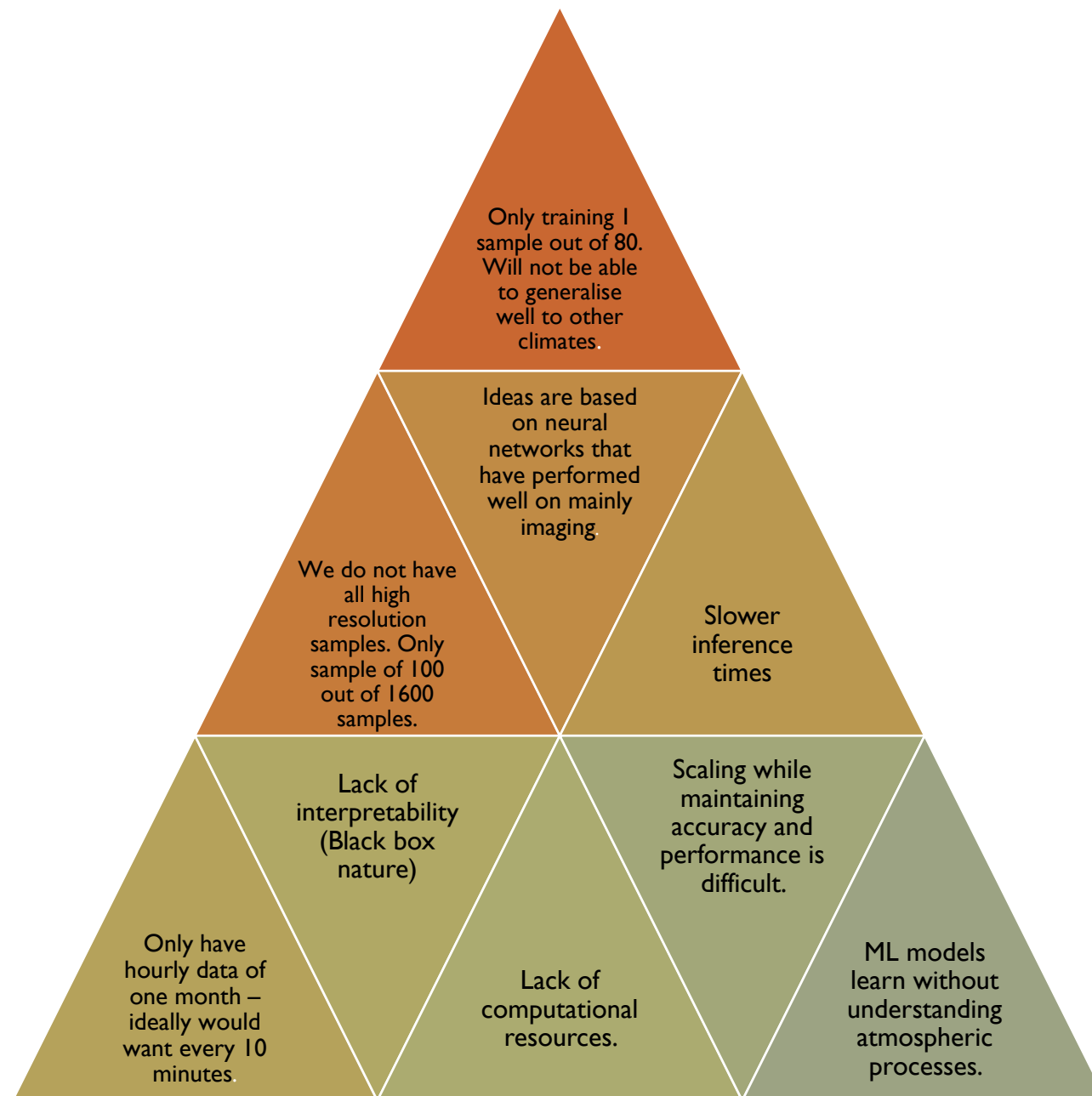
Test our samples to see whether they can still be used for modelling and parametrisation schemes still work.

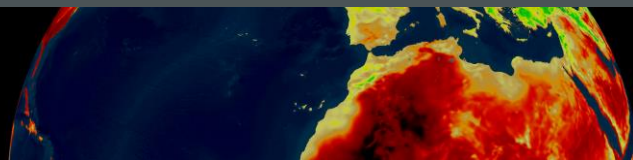
Using data for 2 weeks only so that we can use a standard computer, with modern multi-core processor & adequate memory. Otherwise may need to use GPUs or HPCs, which are commonly used to train neural networks on high dimensional data.





LIMITATIONS





PROBLEMS AND RISKS & ETHICAL ASPECTS

Dependence on AI

Atmospheric models
are sensitive.

Inaccurate results may
misinform decision
making processes

Costs to upscale

System Failures

Constant
updating required

EU AI Act

Accountability issues

Models might develop
biases based on
training data