

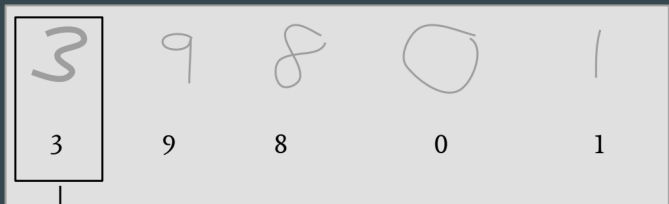
How does stable diffusion work?



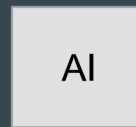
Jax Bulbrook

Digit drawing - Training phase

Training Dataset- a bunch of images along with what they are



Repeat until the loss function is small



Update the AI based on the gradients of that

Compare to original noise

$$\frac{\text{Pred - actual}}{\text{totalInputs}} = \frac{\sum (\hat{n} - n)^2}{\text{totalInputs}}$$

Mean squared error

Noise generator



+

3

Noisy image



=

3

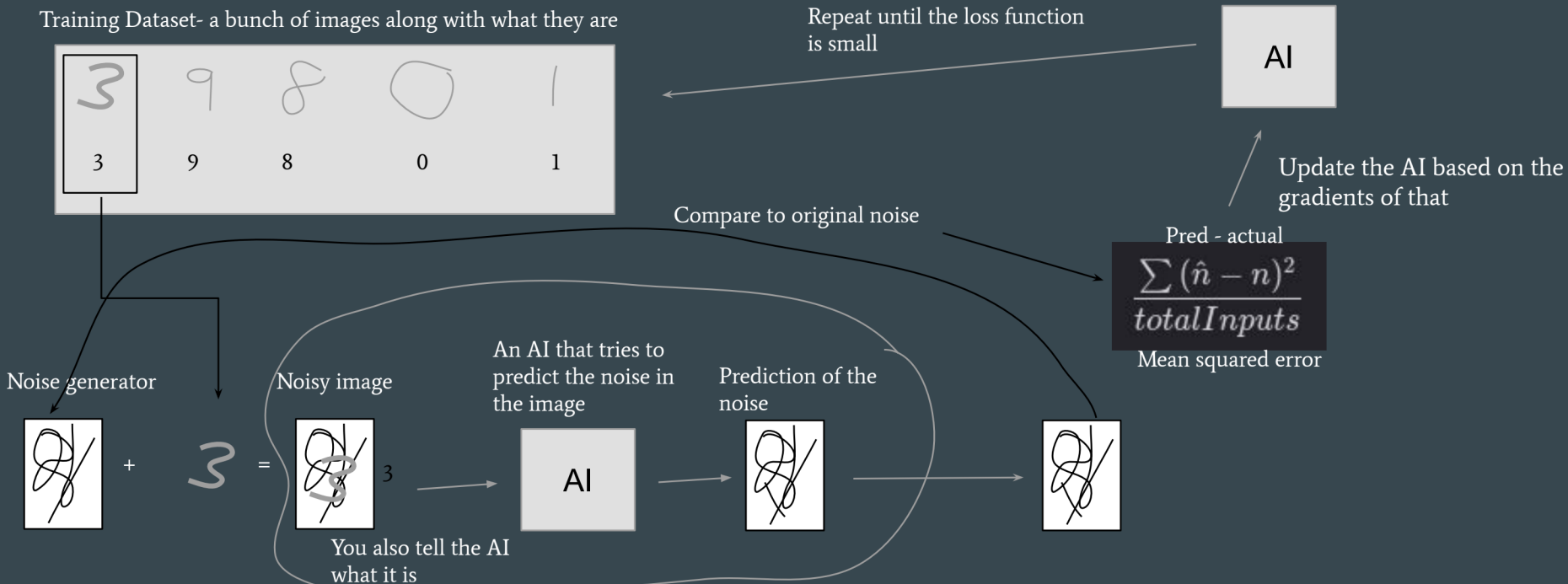
An AI that tries to predict the noise in the image

AI

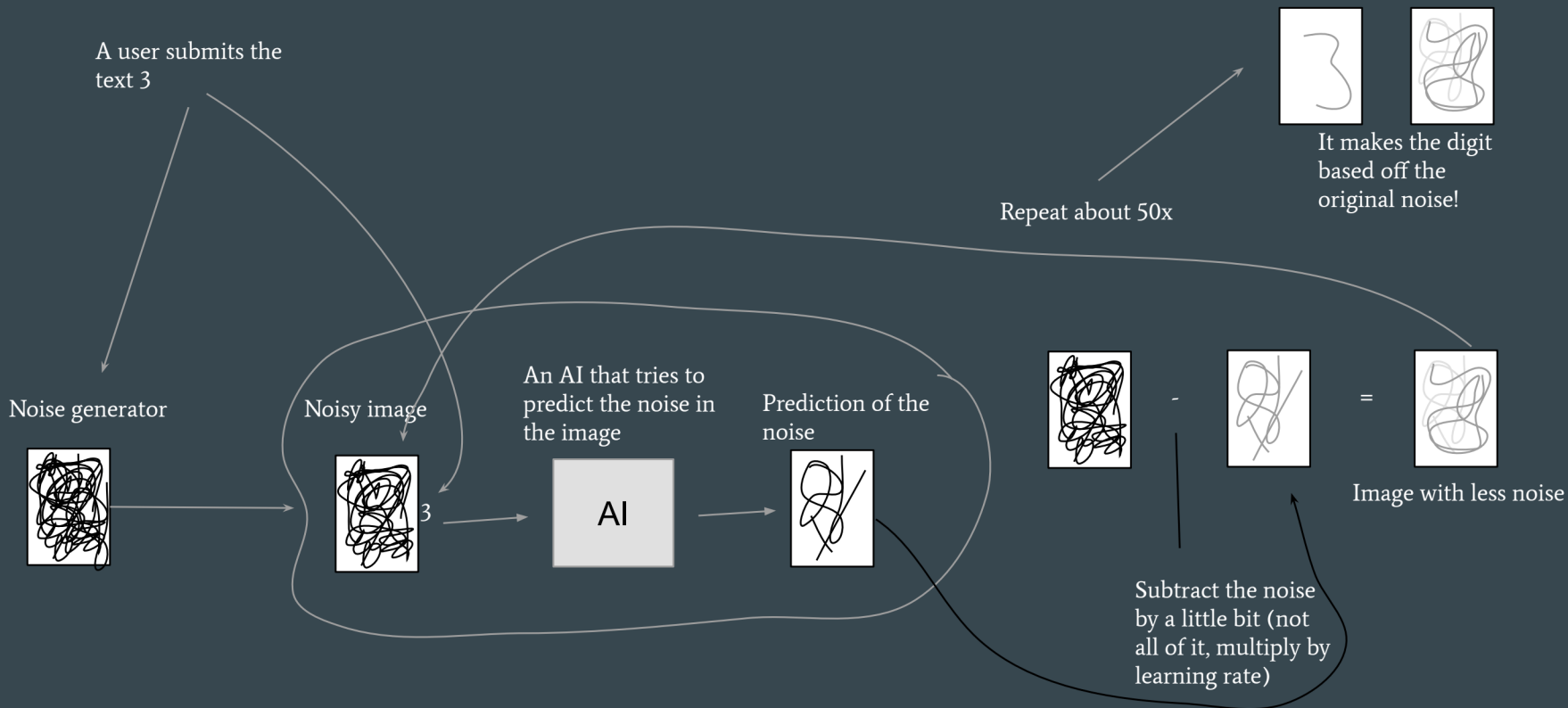
Prediction of the noise



You also tell the AI what it is

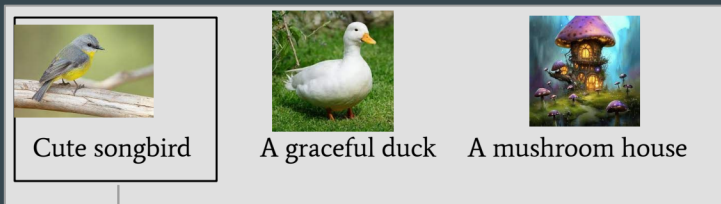


Digit drawing - Production phase



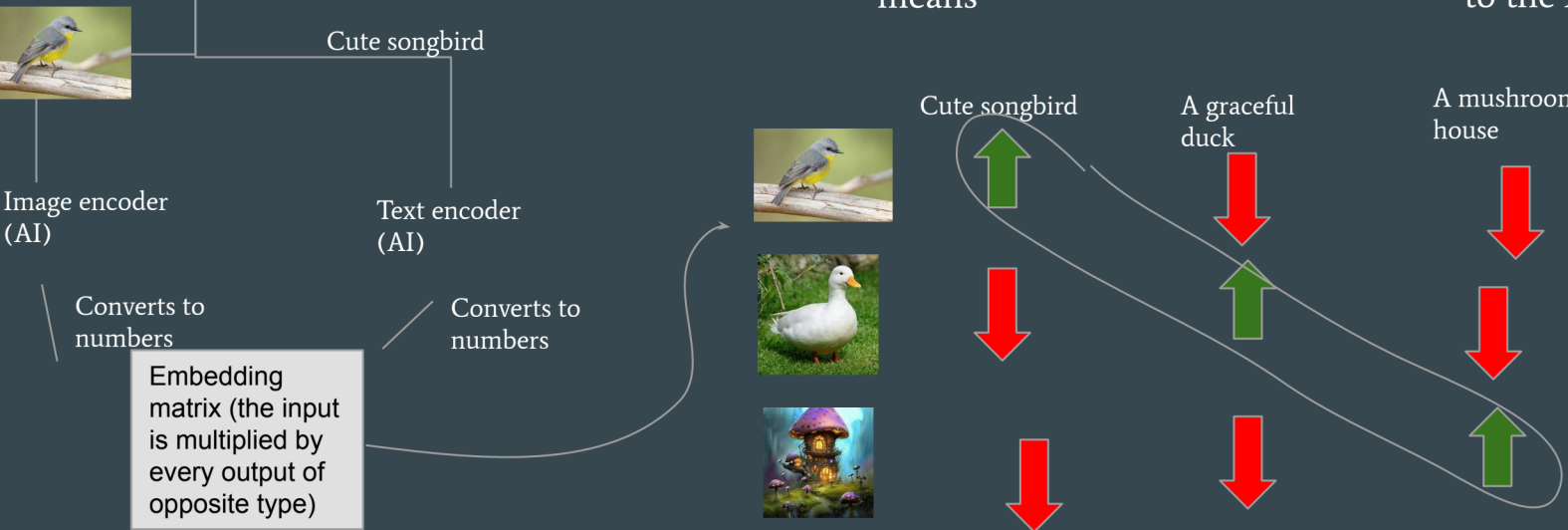
Stable diffusion - understanding phase (CLIP encoder)

Training Dataset- a bunch of images along with what they are, programmatically from alt tag of images found online



Repeat millions of times to train the CLIP encoder - results in the model understanding what the text actually means

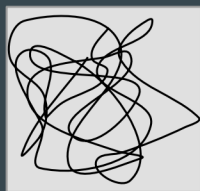
Once trained, the user can pass in a text tag and it will give a list of numbers (embedding) that looks like what the image should look like to the AI (next slide)



Stable Diffusion - Image generation phase (UNET)

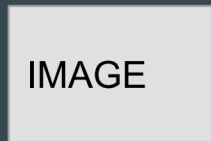
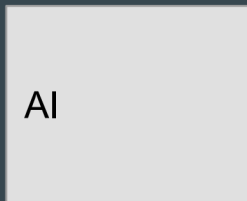
The UNET is trained just like the digit creator, but instead of using a description of “3” or “9”, it would be the list of numbers

List of numbers (embedding) that should look like an image, came from user



Random noise generator

01010010010101.....
embedding



Predicts where noise is, subtract that times learning rate, repeat, etc x50