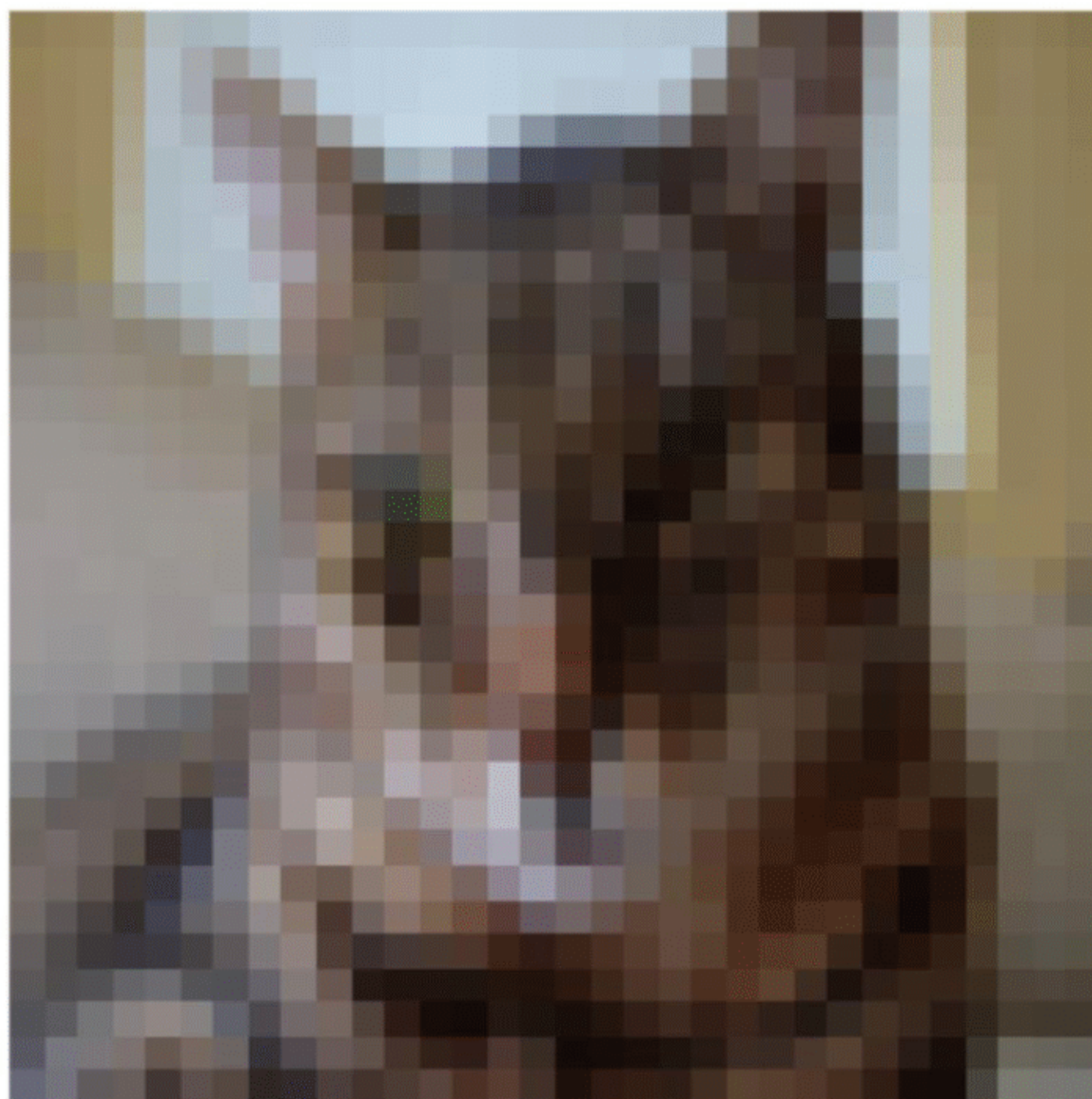
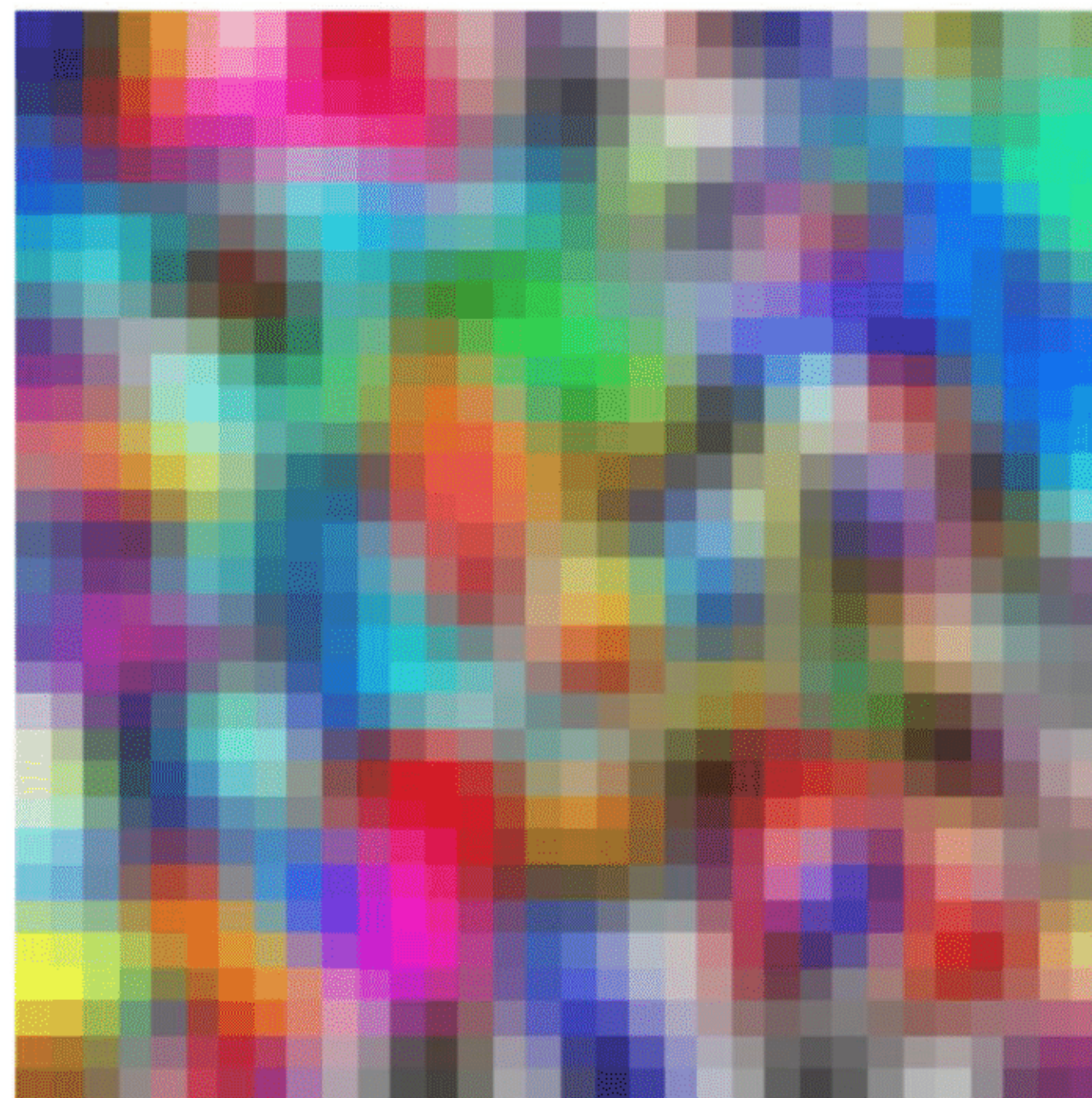


Original cat



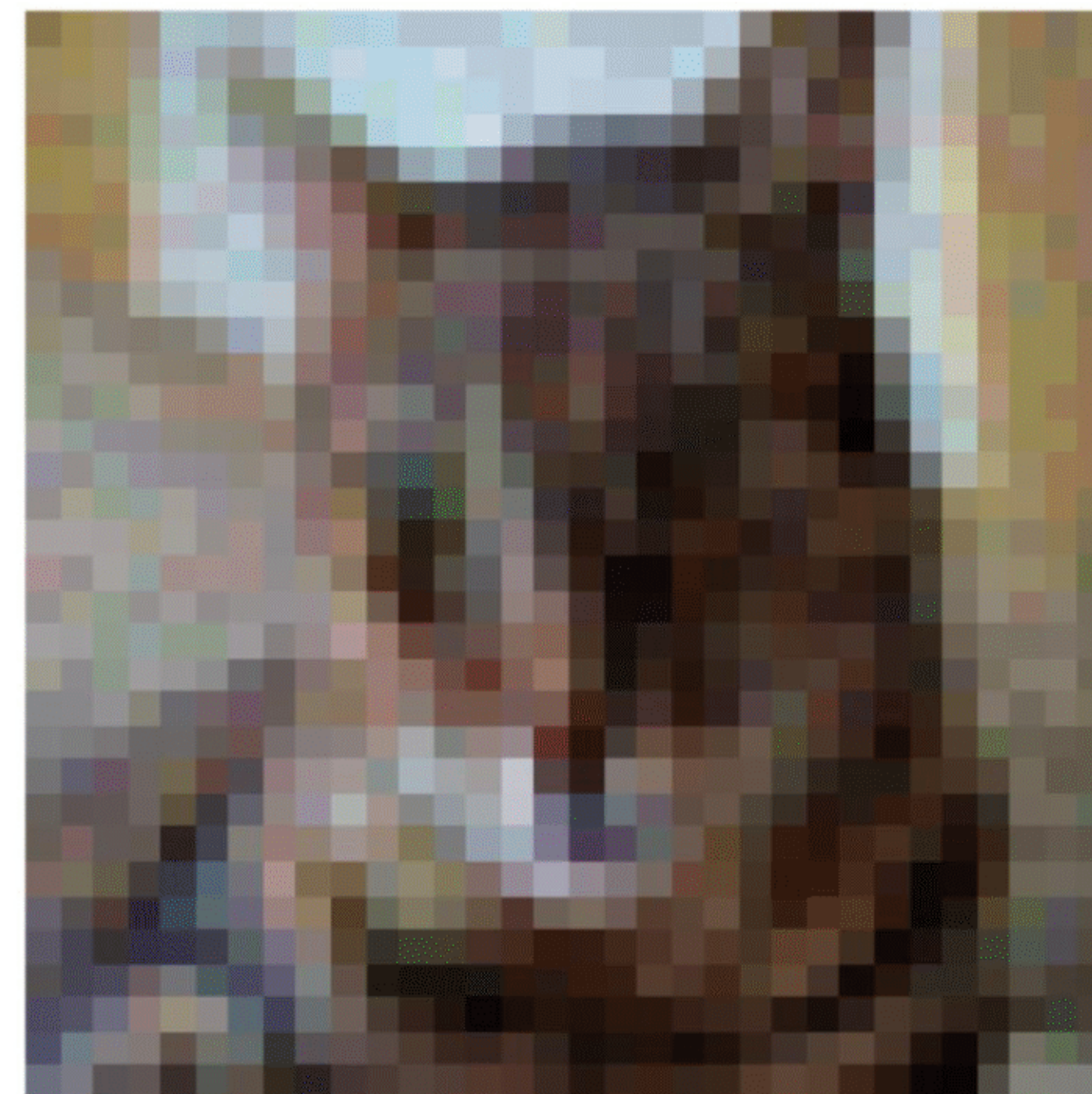
Perturbation



+



Adversarial sample



x

$$0.031 \times \text{sign}(\nabla_x \ell(g(x; \theta), y))$$

$$x + \text{sign}(\nabla_x \ell(g(x; \theta), y))$$

“cat”

94.6% confidence

“dog”

99.8% confidence