

# ClearStamp

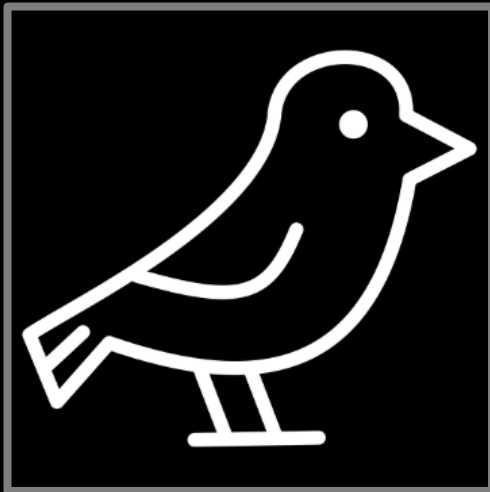
A Human-Visible and Robust Model-Ownership Proof  
based on Transposed Model Training

Torsten Krauß, Jasper Stang, and Alexandra Dmitrienko

University of Würzburg, Germany

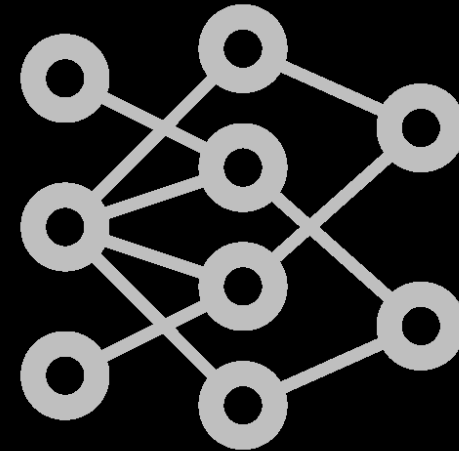
33<sup>rd</sup> USENIX Security Symposium

*Image*



**WATERMARKED**

*ML Model*



**WATERMARKED**

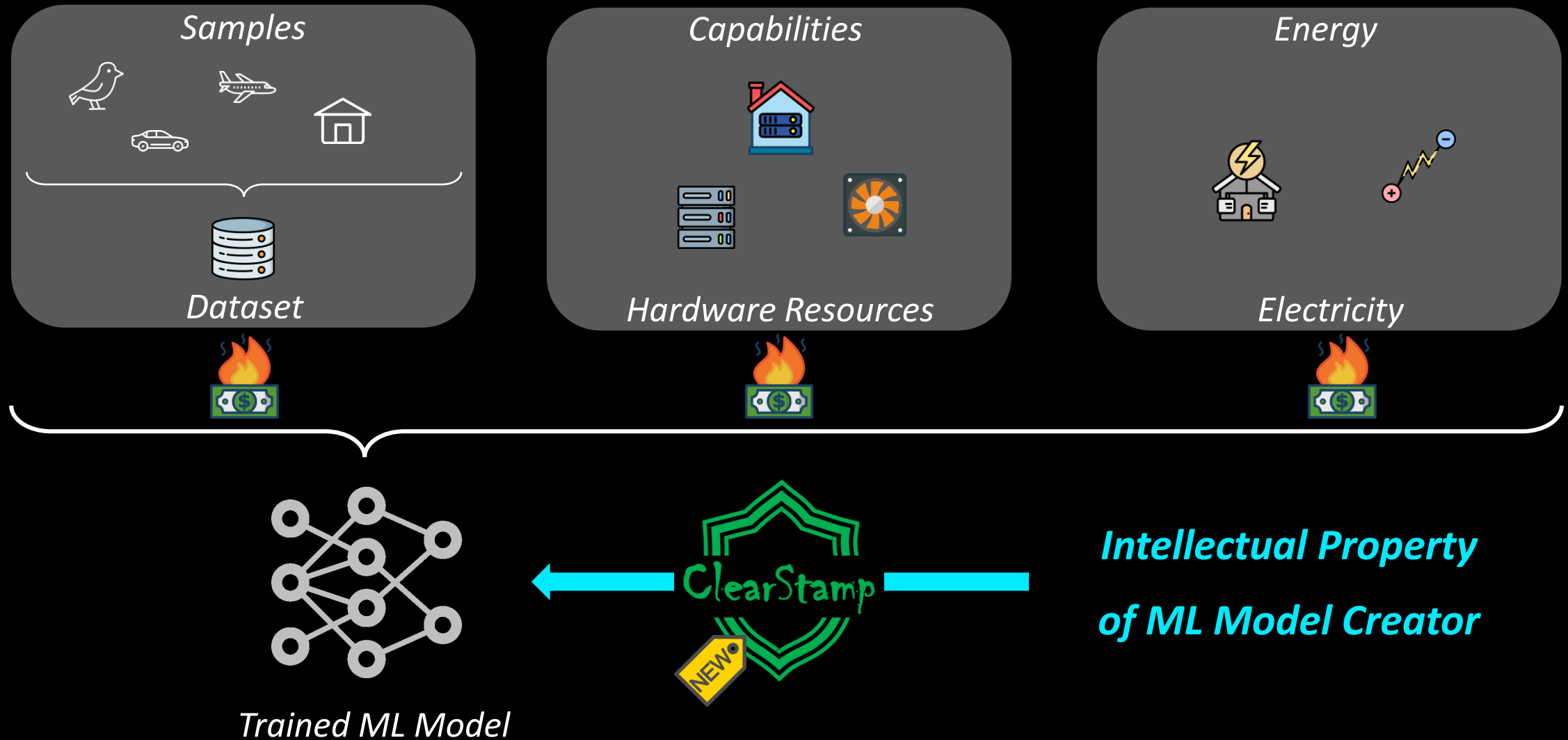
*Image*



*ML Model*

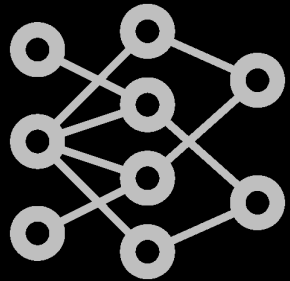


# ML Models are Intellectual Property



# Intellectual Property Violations

*Trained ML Model*



- *Theft*



- *Illicit Utilization*



- *Unauthorized Resell*



- *Unwarranted Modifications*



# Intellectual Property Protection

*Image*

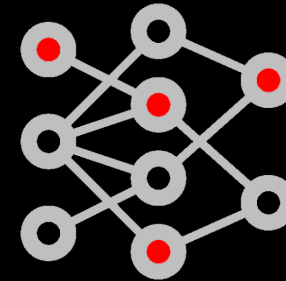


*Add a Visible Stamp*

*Same Principle*

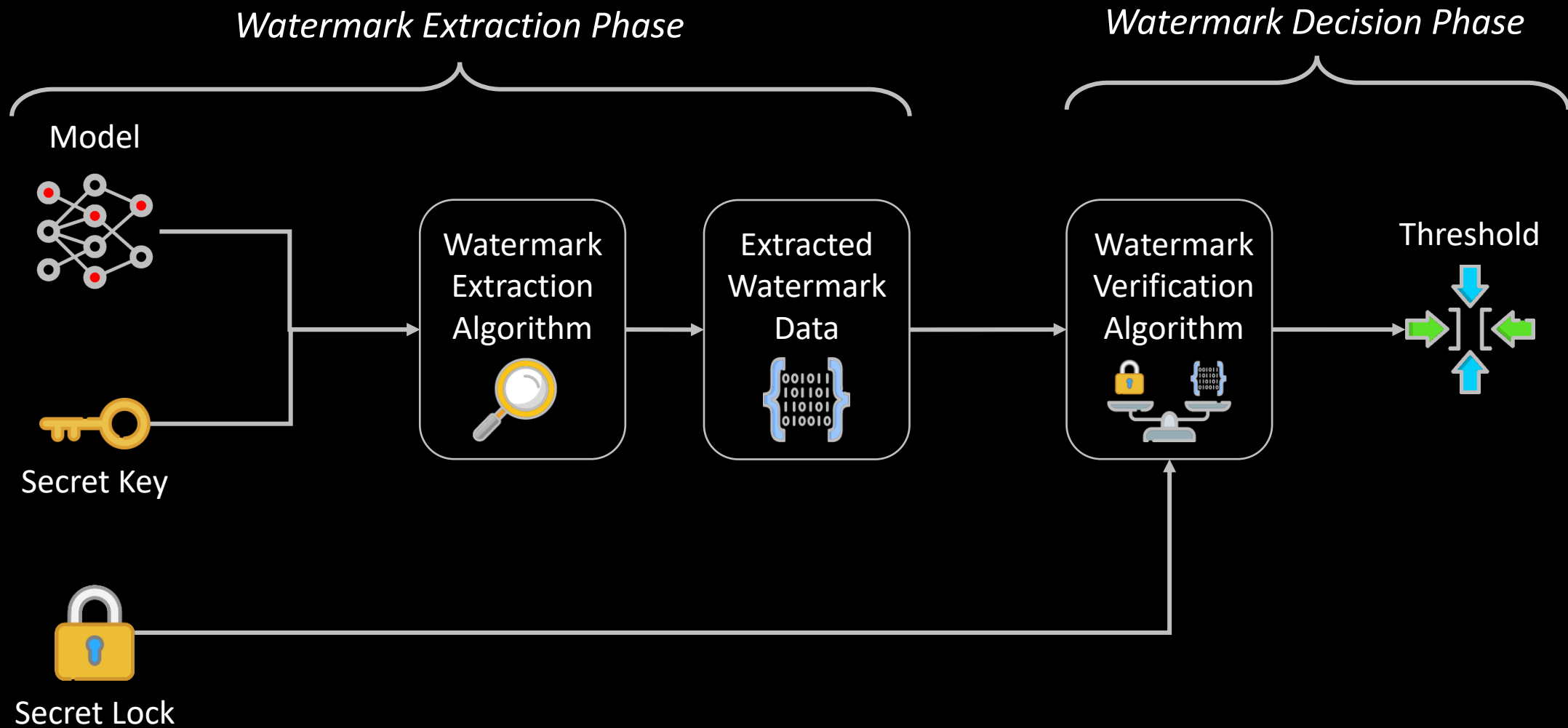


*Trained ML Model*

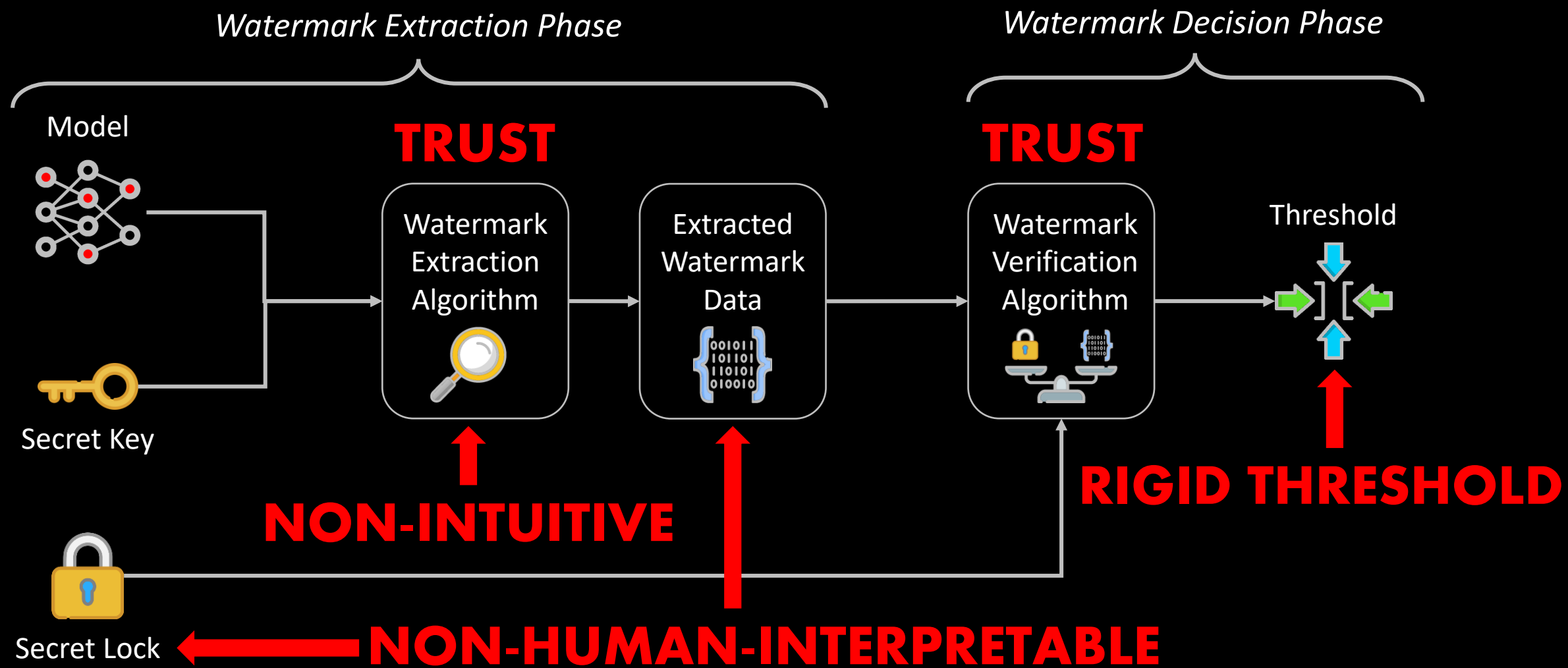


*Hide Unique Information in Parameters*

# ML Model Watermark Verification

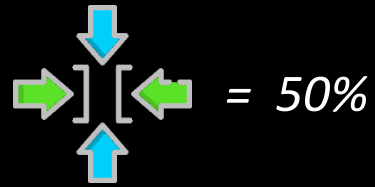


# ML Model Watermark Verification





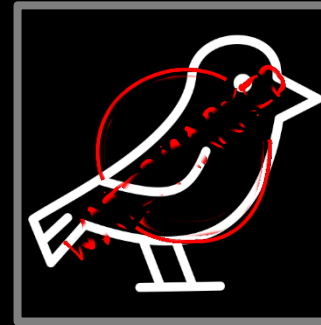
# Rigid Threshold VS. Partially Removed Watermarks



100 % Watermark



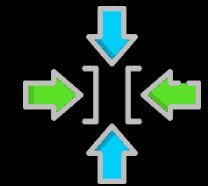
85 % Watermark



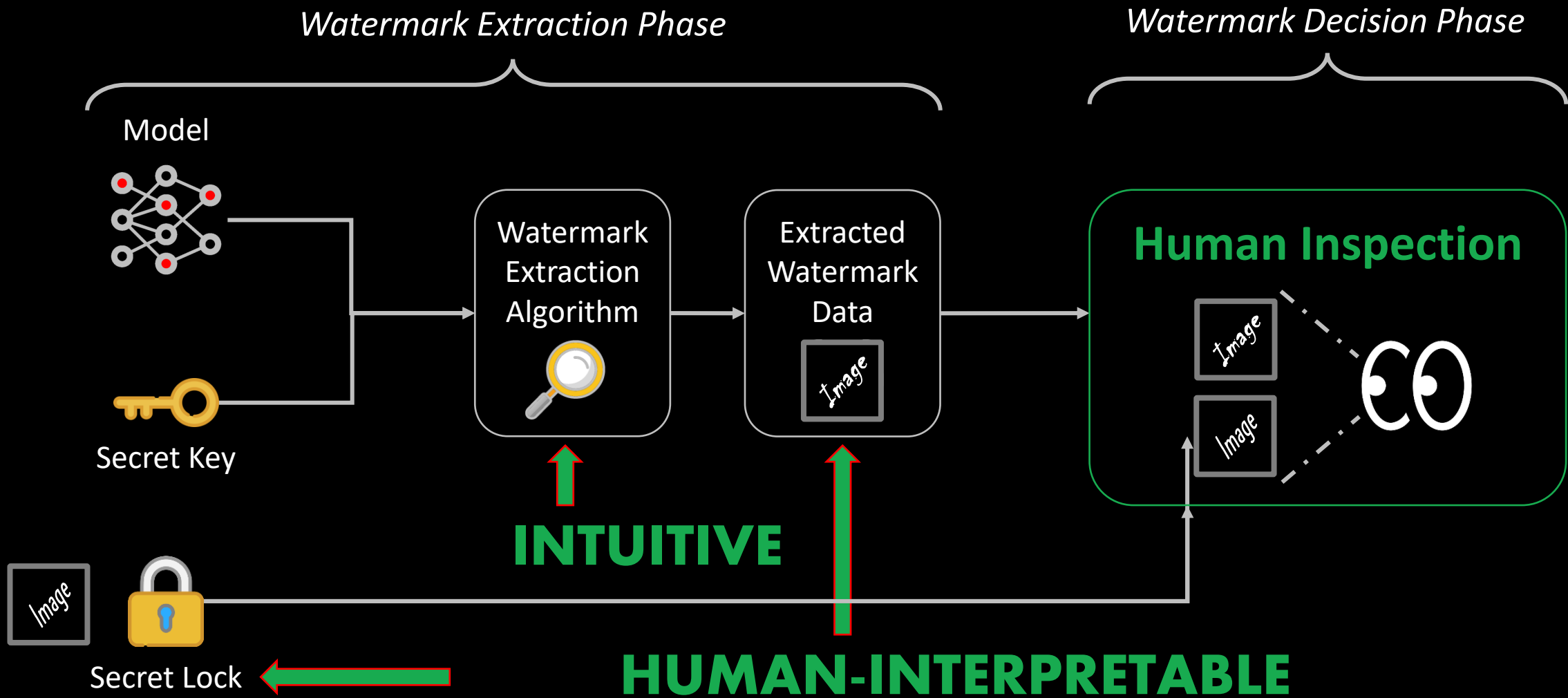
49 % Watermark



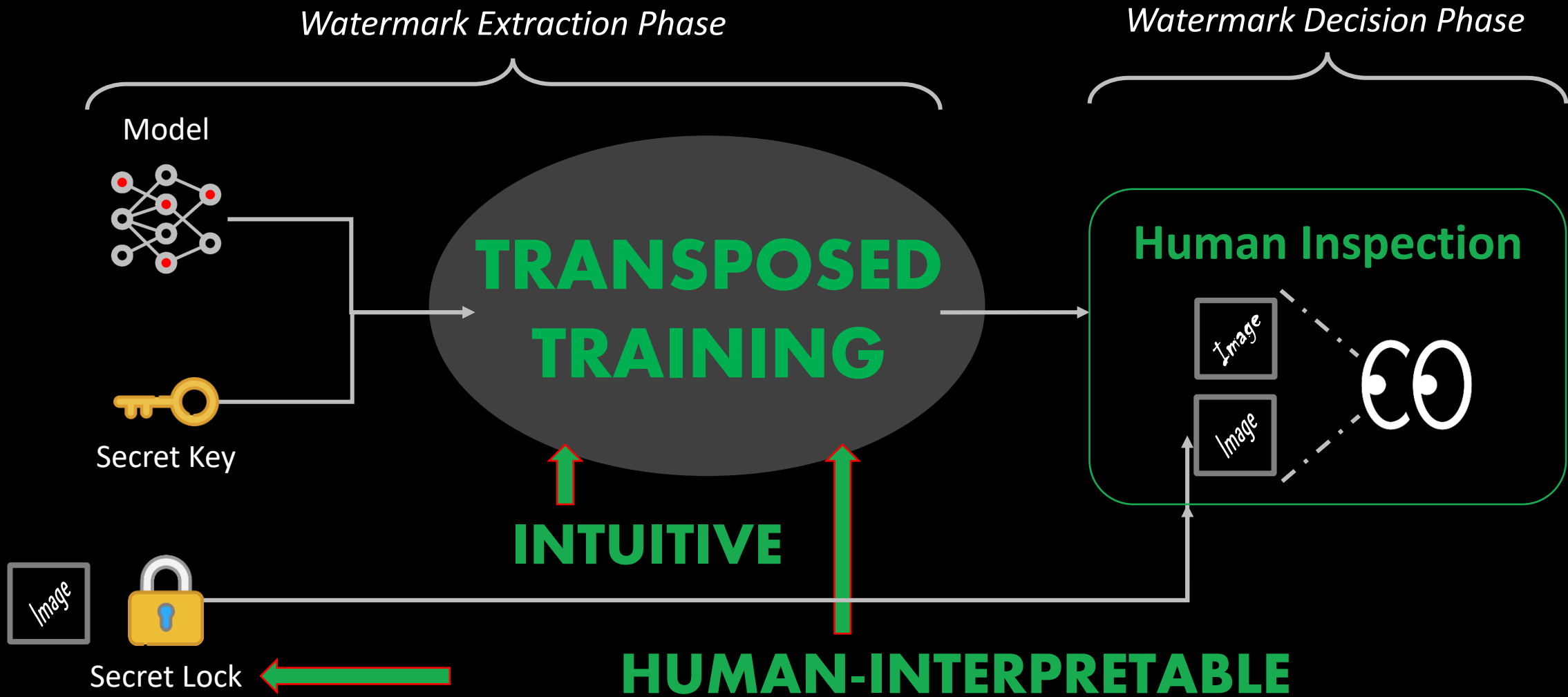
20 % orange Watermark



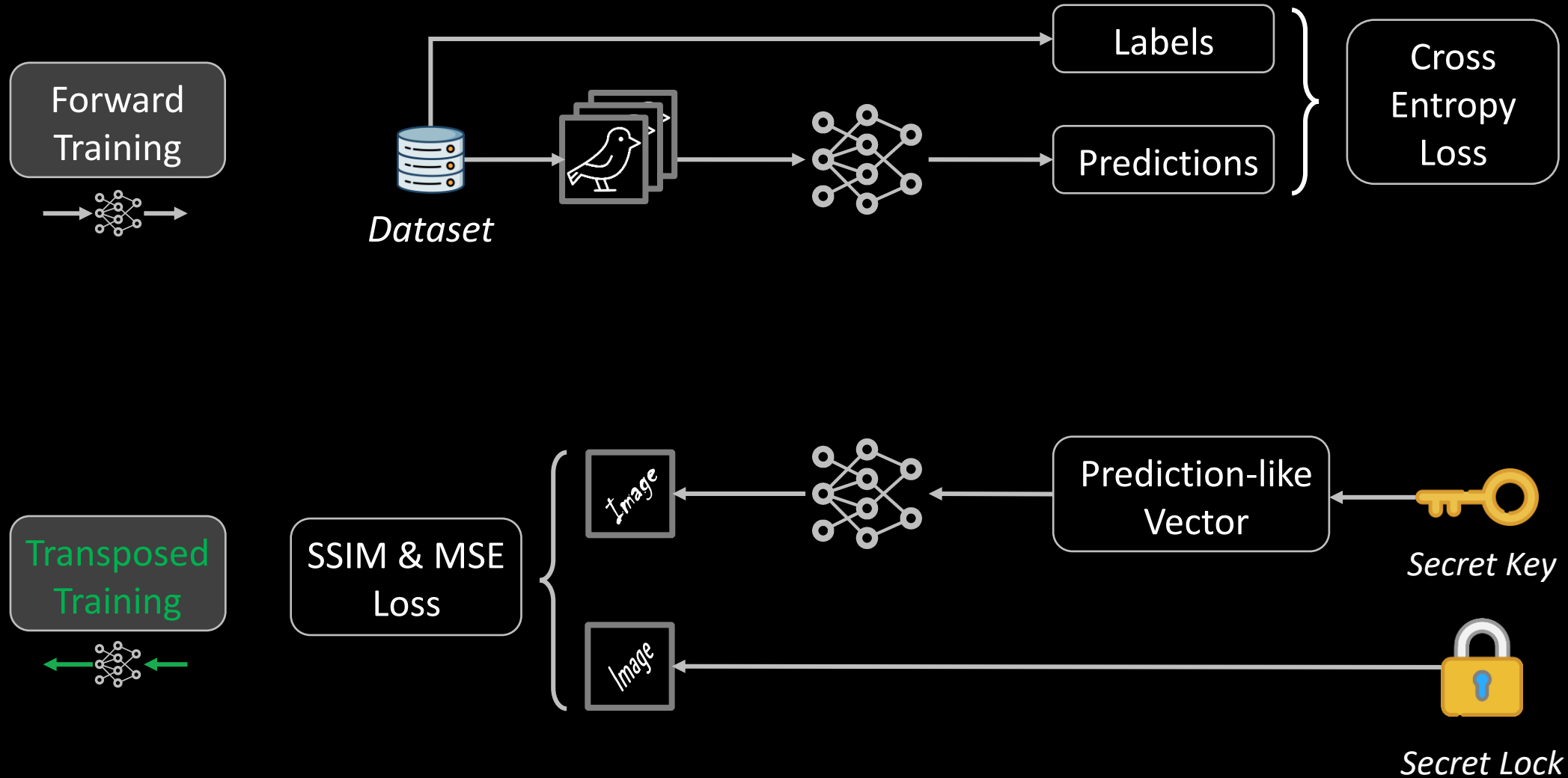
# ClearStamp - Principle



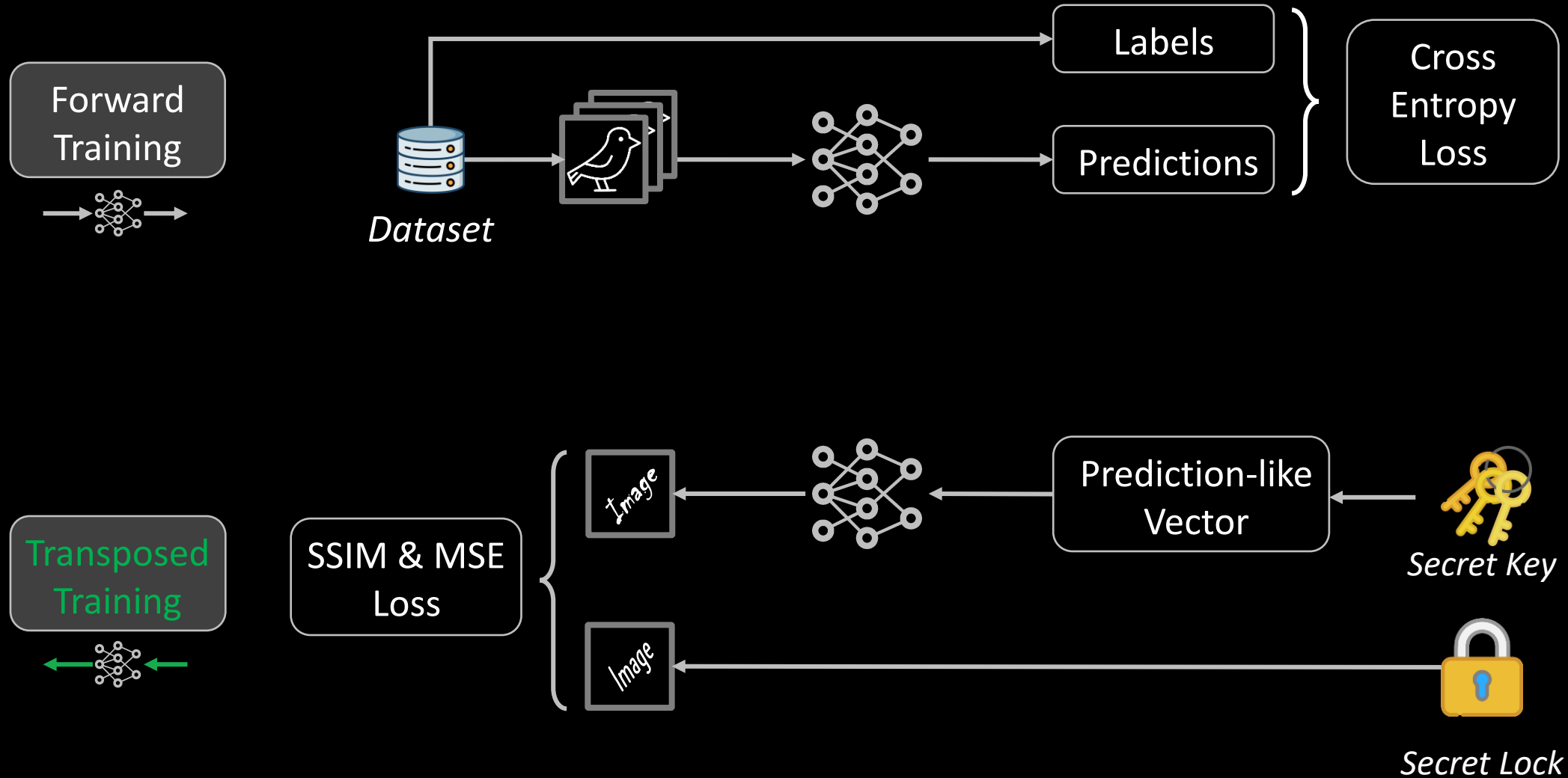
# ClearStamp - Principle



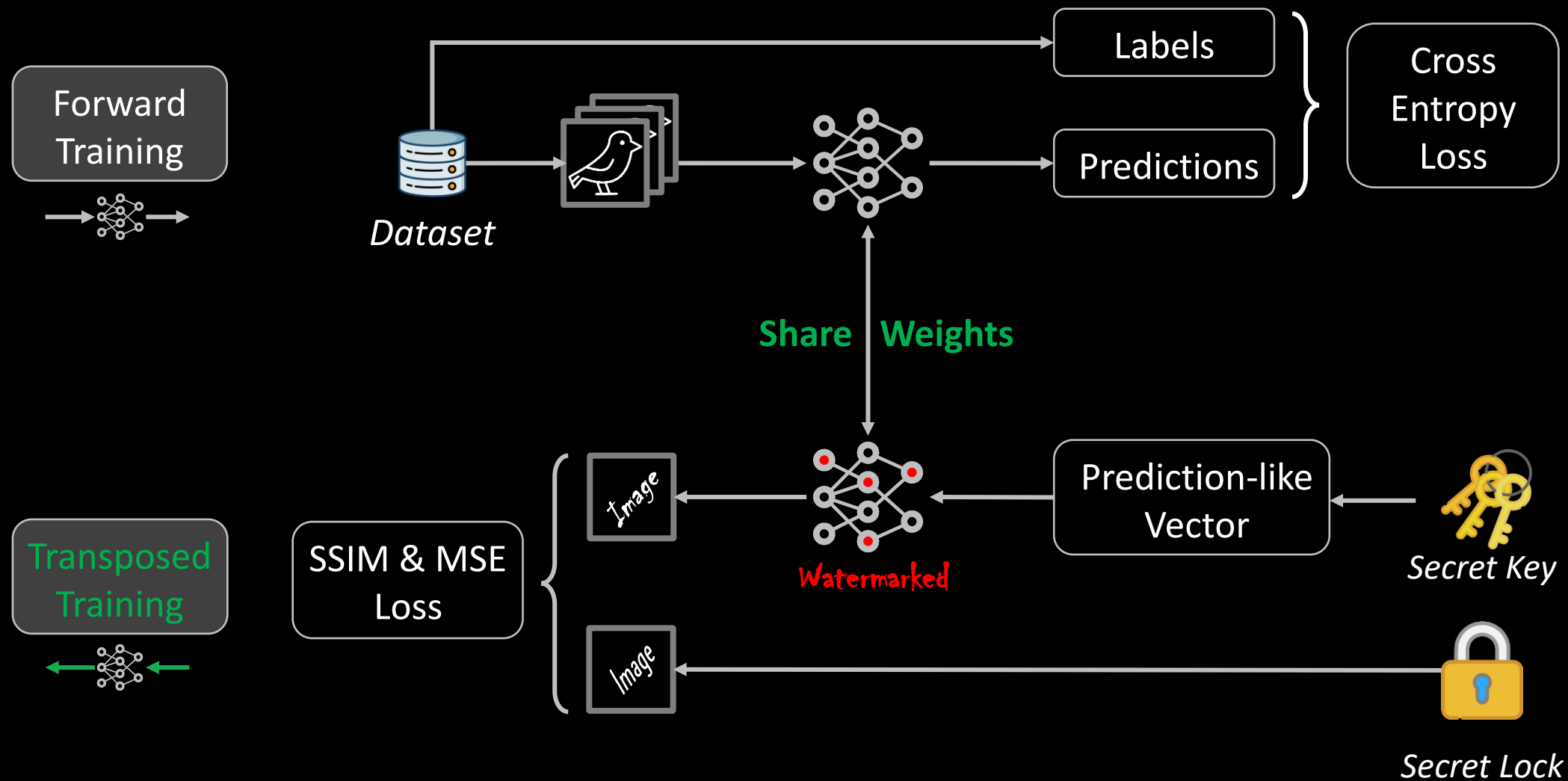
# Transposed Training for Watermarking



# Transposed Training for Watermarking



# Transposed Training for Watermarking



# Transposed Training - Details

Model Component	Forward Model	Transposed Model
Linear Layer	$y = x \cdot w^T + b$	$x = (y - b) \cdot w$
Batch Normalization	$y = \frac{x - E(x)}{\sqrt{Var(x) + \varepsilon}} \cdot \gamma + \beta$	$x = \frac{(y - \beta) \cdot \sqrt{1 + \varepsilon}}{\gamma}$ with $E(x) = 0, Var(x) = 1$
Convolutions [1]	Replace with deconvolutions [2]	
Pooling Layer [3]	Replace with Interpolations [4, 5]	
Dropout Layers [6]	Keep same dropout	
Activation Functions	Use same activation, e.g., ReLU [7]	
Skip Connections	Fixate skip connections	

[1] Y. Lecun, L. Bottou, Y. Bengio, and P. Haffner. Gradient-based learning applied to document recognition. Proceedings of the IEEE, 1998.

[2] Matthew D Zeiler, Dilip Krishnan, Graham W Taylor, and Rob Fergus. Deconvolutional networks. CVPR, 2010.

[3] Afia Zafar, Muhammad Aamir, Nazri Mohd Nawi, Ali Arshad, Saman Riaz, Abdulrahman Alruban, Ashit Kumar Dutta, and Sultan Almotairi. A comparison of pooling methods for convolutional neural networks. AppliedSciences, 2022.

[4] Olivier Rukundo and Hanqiang Cao. Nearest Neighbor Value Interpolation. IJACSA, 2012.

[5] Olivier Rukundo and Bodhaswar T Maharaj. Optimization of Image Interpolation based on Nearest Neighbour Algorithm. VISAPP, 2014.

[6] Nitish Srivastava, Geoffrey Hinton, Alex Krizhevsky, Ilya Sutskever, and Ruslan Salakhutdinov. Dropout: A simple way to prevent neural networks from overfitting. Journal of Machine Learning Research, 2014.

[7] Abien Fred Agarap. Deep Learning using Rectified Linear Units (ReLU). arXiv preprint arXiv:1803.08375, 2018.

# Transposed Training - Details

Model Component	Forward Model	Transposed Model
Linear Layer	$y = x \cdot w^T + b$	$x = (y - b) \cdot w$
Batch Normalization	$y = \frac{x - E(x)}{\sqrt{Var(x) + \epsilon}} \cdot \gamma + \beta$	$x = \frac{(y - \beta) \cdot \sqrt{1 + \epsilon}}{\gamma}$ with $E(x) = 0, Var(x) = 1$
Convolutions [1]		Replace with deconvolutions [2]
Pooling Layer [3]		Replace with Interpolation [4, 5]
Dropout Layers [6]		Keep same dropout
Activation Functions		Use same activation, e.g., ReLU [7]
Skip Connections		Fixate skip connections

Details in Paper

[1] Y. Lecun, L. Bottou, Y. Bengio, and P. Haffner. Gradient-based learning applied to document recognition. Proceedings of the IEEE, 1998.

[2] Matthew D Zeiler, Dilip Krishnan, Graham W Taylor, and Rob Fergus. Deconvolutional networks. CVPR, 2010.

[3] Afia Zafar, Muhammad Aamir, Nazri Mohd Nawi, Ali Arshad, Saman Riaz, Abdulrahman Alruban, Ashit Kumar Dutta, and Sultan Almotairi. A comparison of pooling methods for convolutional neural networks. AppliedSciences, 2022.

[4] Olivier Rukundo and Hanqiang Cao. Nearest Neighbor Value Interpolation. IJACSA, 2012.

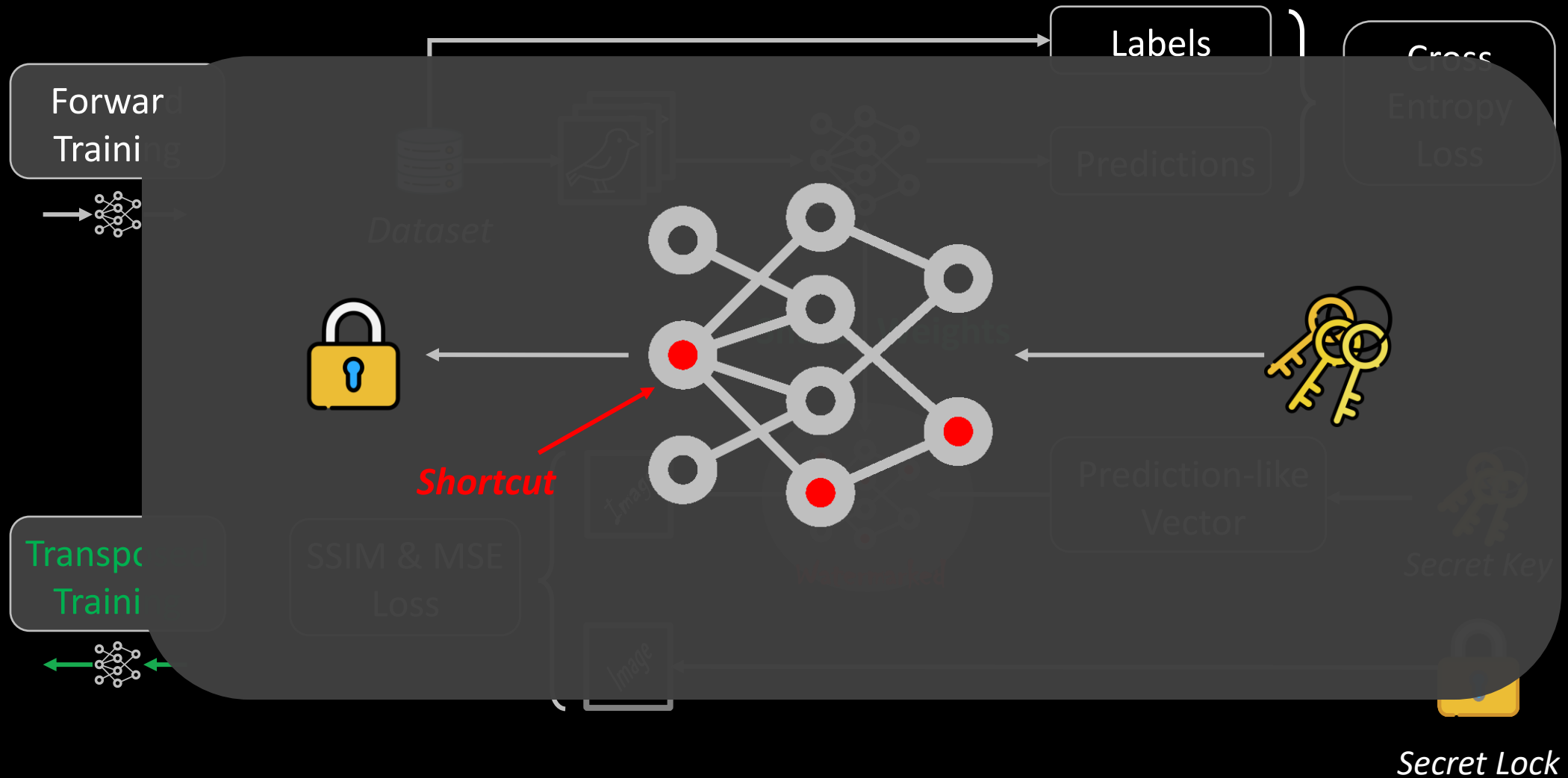
[5] Olivier Rukundo and Bodhaswar T Maharaj. Optimization of Image Interpolation based on Nearest Neighbour Algorithm. VISAPP, 2014.

[6] Nitish Srivastava, Geoffrey Hinton, Alex Krizhevsky, Ilya Sutskever, and Ruslan Salakhutdinov. Dropout: A simple way to prevent neural networks from overfitting. Journal of Machine Learning Research, 2014.

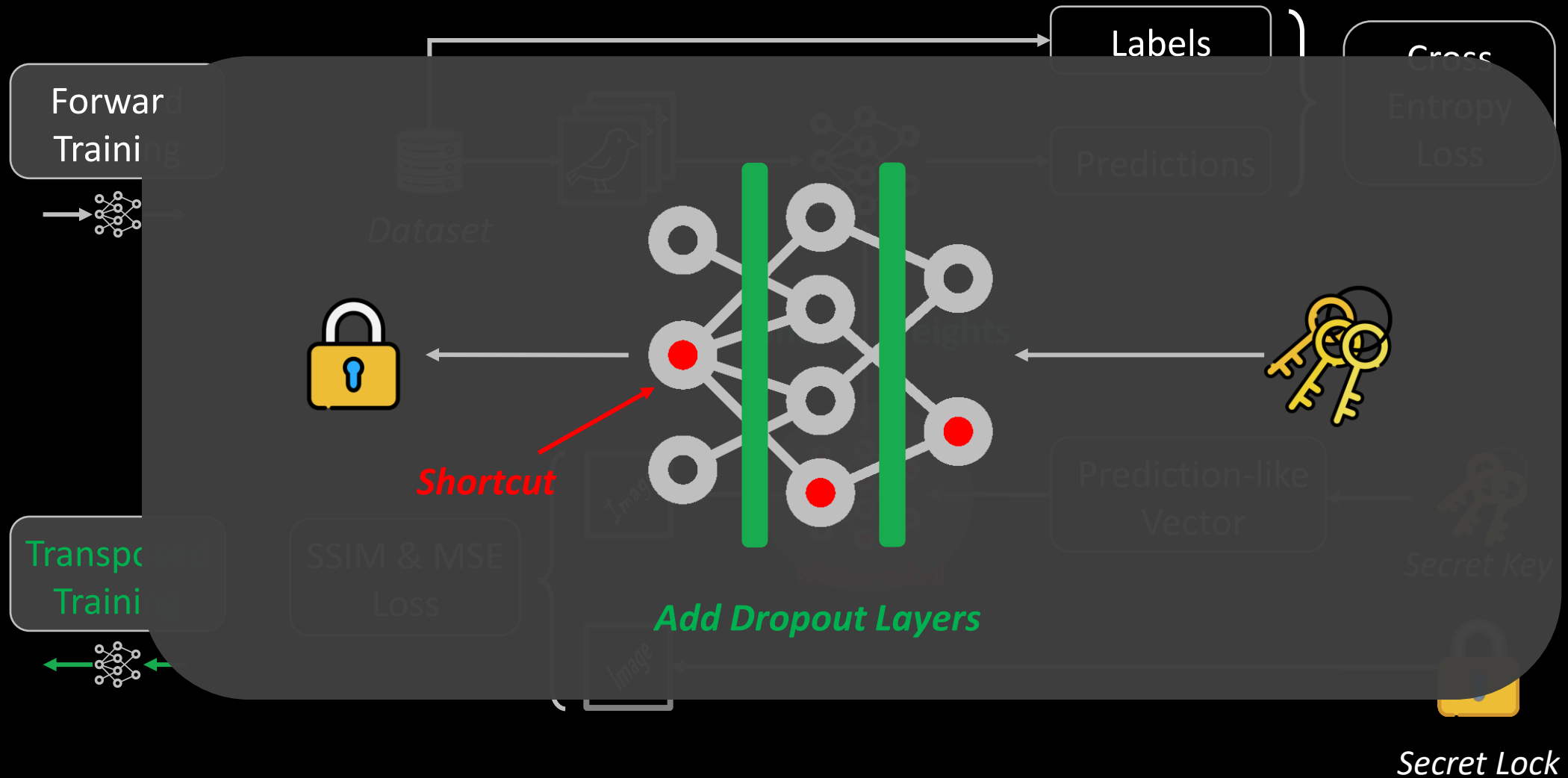
[7] Abien Fred Agarap. Deep Learning using Rectified Linear Units (ReLU). arXiv preprint arXiv:1803.08375, 2018.



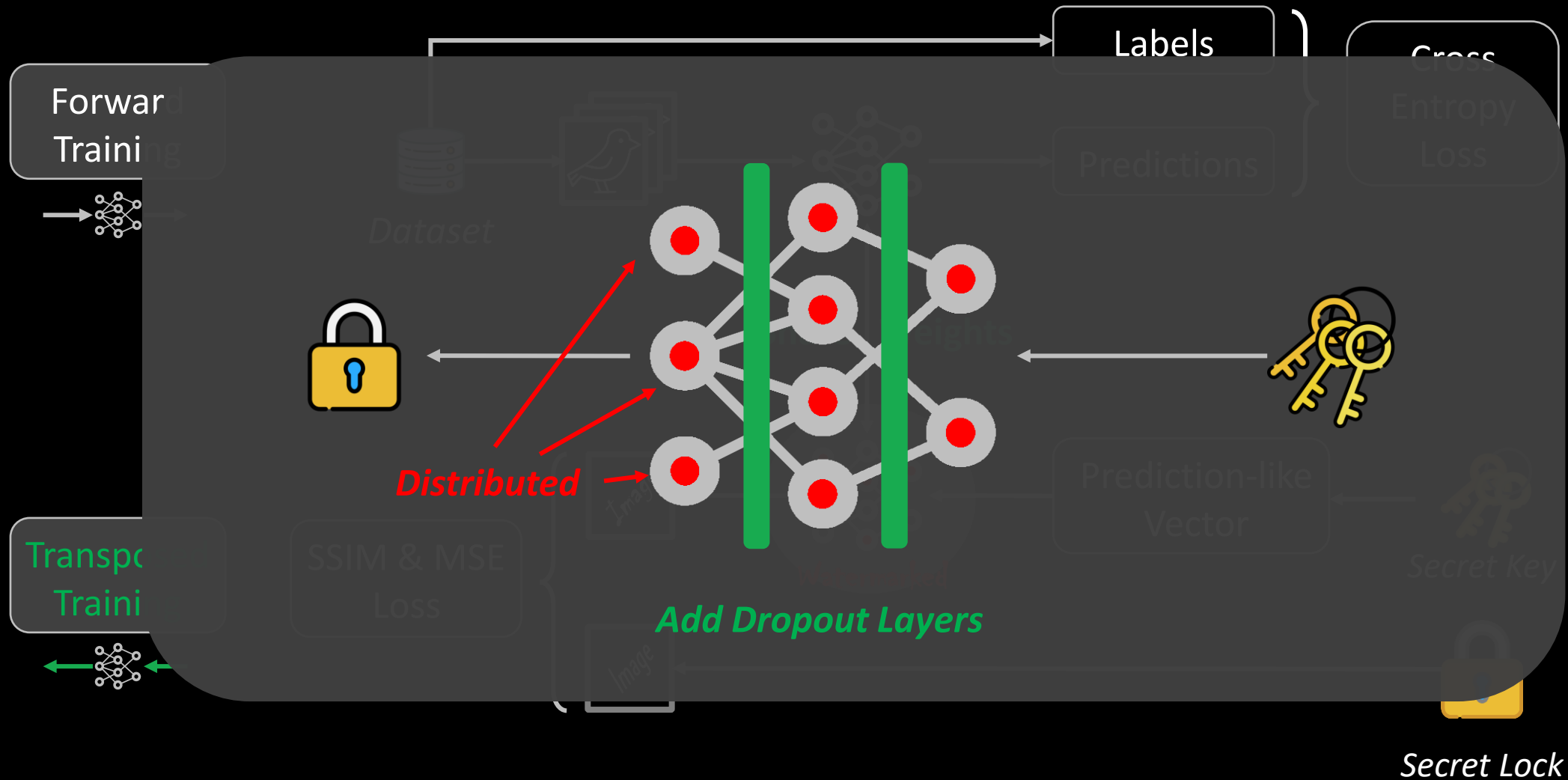
# Parameter Entanglement



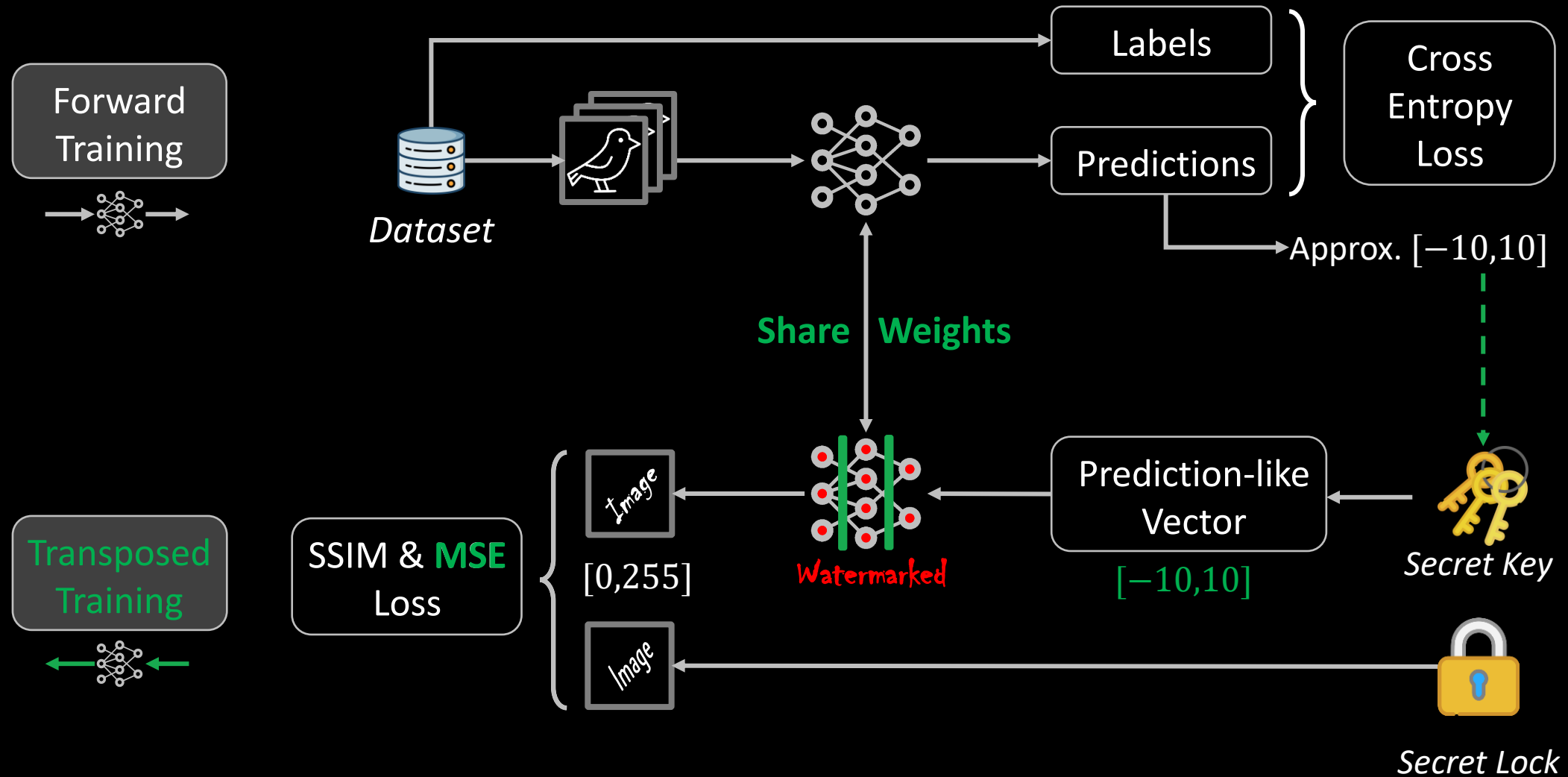
# Parameter Entanglement



# Parameter Entanglement



# Parameter Entanglement



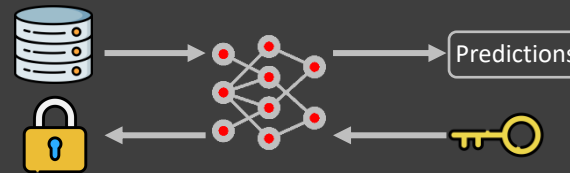
# ClearStamp - Workflow



1 Watermark Hardening

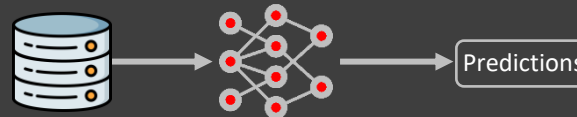


2 Constraint Training



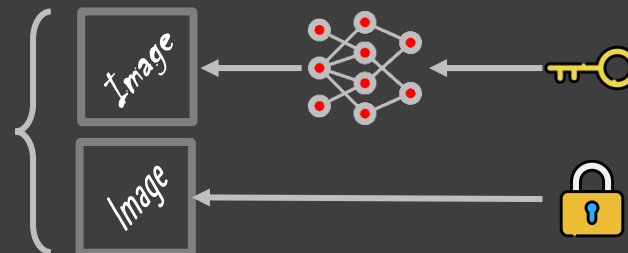
Legal / Illegal Model Distribution

3 3<sup>rd</sup> Party Manipulation 

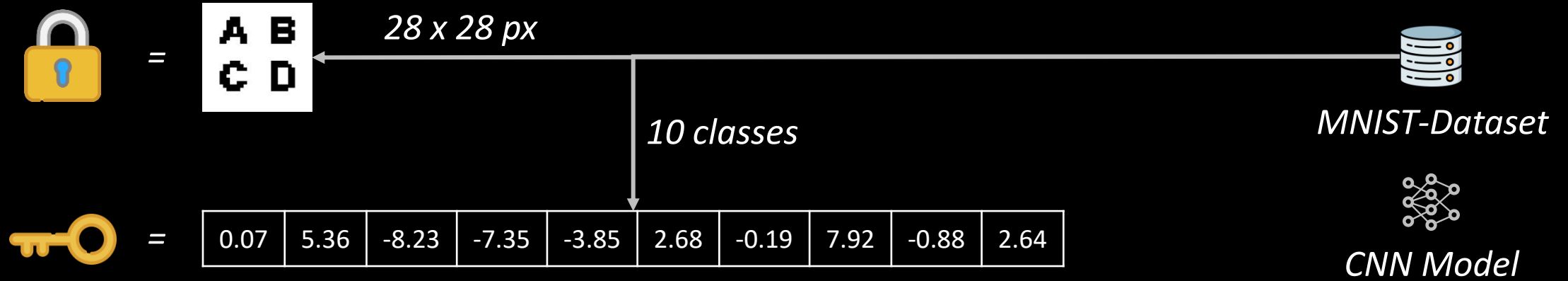


Copyright Infringement?

4 Watermark Testing



# ClearStamp - Evaluation



<i>Model Performance</i>	10.22 %	89.88 %	8.57 %	88.37 %	87.69 %
<i>Extracted Watermark</i>					

# ClearStamp - Evaluation



=



10 classes



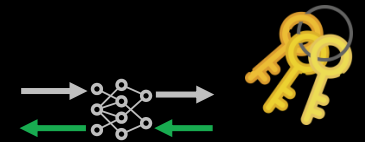
MNIST-Dataset



=



CNN Model



Model Performance

10.22 %

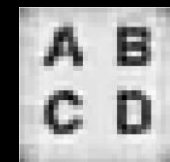
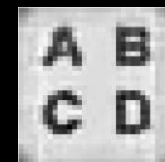
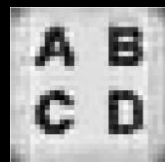
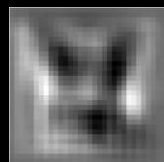
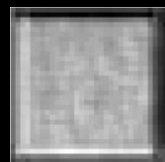
89.88 %

8.57 %

88.57 %

87.69 %

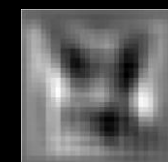
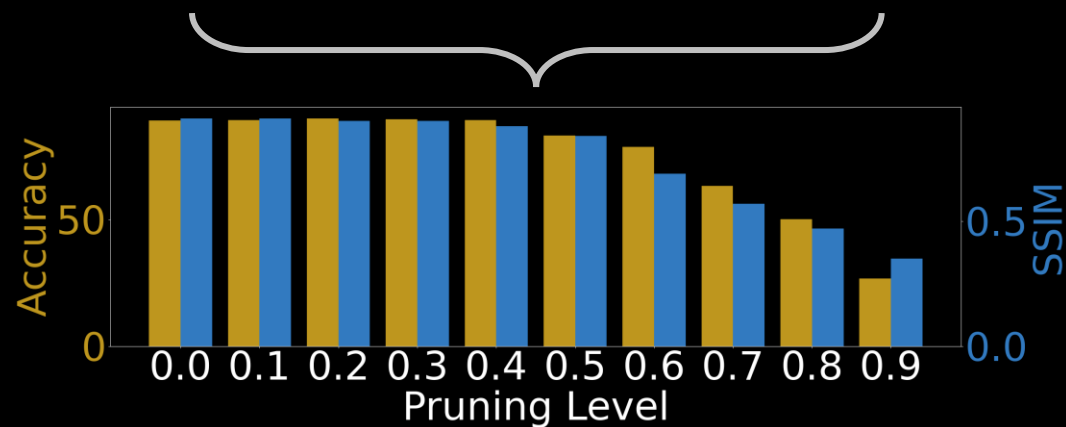
Extracted Watermark



# ClearStamp - Evaluation



	<i>Fine-Tuning</i>			<i>Pruning</i>			<i>Fine-Pruning</i>
	<i>Same LR</i>	$1/_{10}$ LR	$1/_{10}$ LR & <i>unseen data</i>	60 %	80 %	90 %	$1/_{10}$ LR & 40 %
<i>Performance</i>	87.42 %	89.86 %	97.02 %	78.56 %	50.10 %	26.76 %	89.79 %
<i>Extracted Watermark</i>							





# ClearStamp - Evaluation



*Erase Watermark*

*Performance*

*Capacity*

*Runtime*



83.89 %

70.38 %

56.39 %

46.52 %

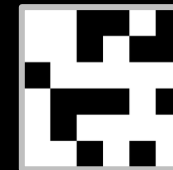
MiT License Text

**8,544 bits**



Bit Error Rate

5.92 %



Dot code

Hardening

53.48 s (**one time**)

Training

67.62 s → 94.39 s

**+ 39.58 %**

Extraction

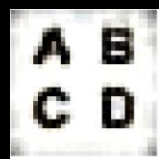
0.02 s

# ClearStamp - Evaluation

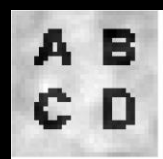
*CNN*  
*CIFAR-10*



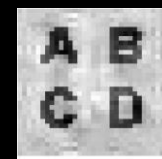
*CNN*  
*GTSRB*



*Small CNN*  
*MNIST*



*CNN + Batchnorm*  
*MNIST*



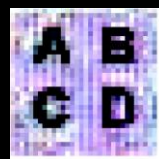
*ResNet-18*  
*CIFAR-10*



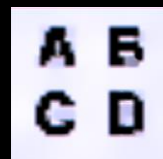
*ResNet-18*  
*GTSRB*



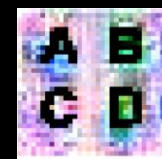
*ResNet-34*  
*CIFAR-10*



*VGG11*  
*CIFAR-100*



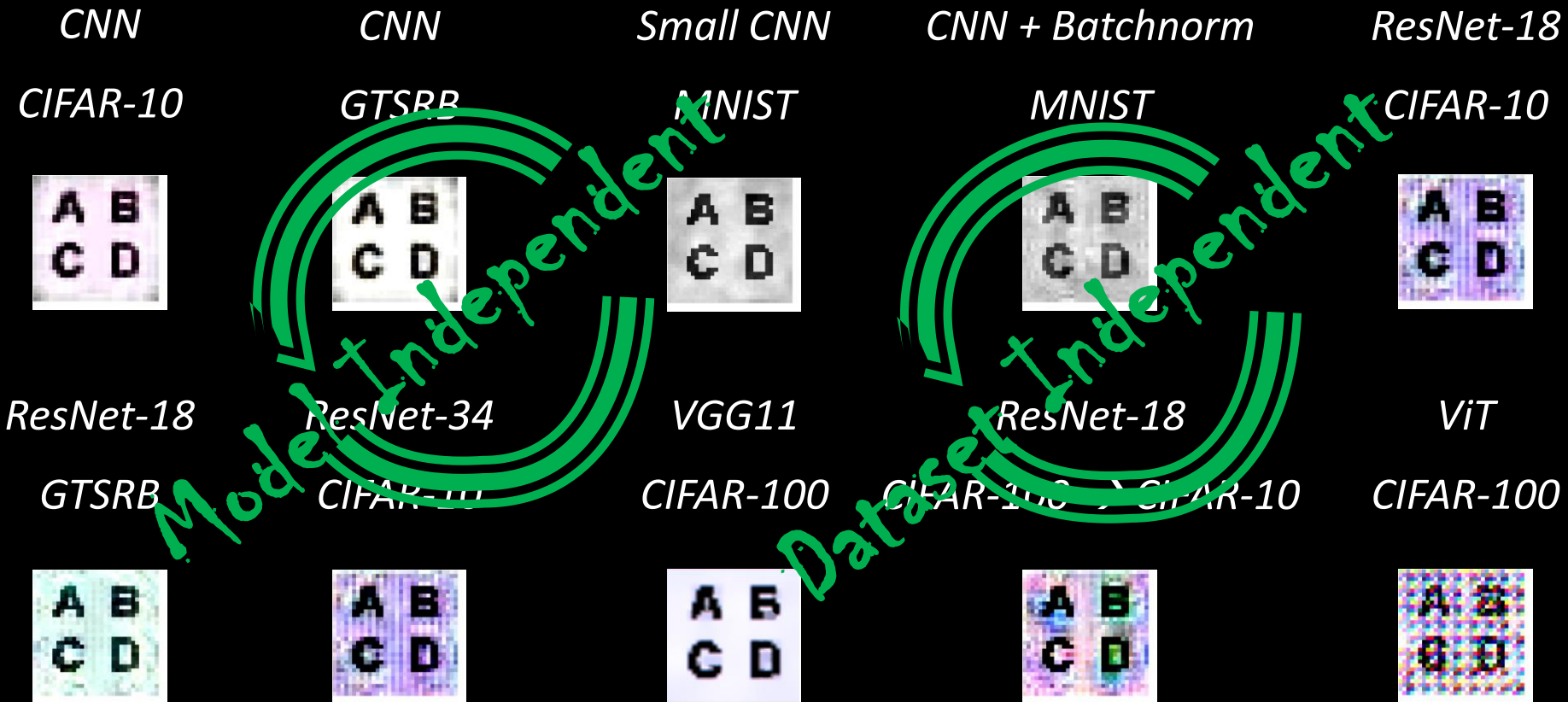
*ResNet-18*  
*CIFAR-100* → *CIFAR-10*



*ViT*  
*CIFAR-100*



# ClearStamp - Evaluation



# Conclusion



Copyright Infringement of ML models



Watermarking of ML models



- Non-intuitive algorithms
- Non-human-interpretable
- Rigid threshold
- Partially removed watermarks

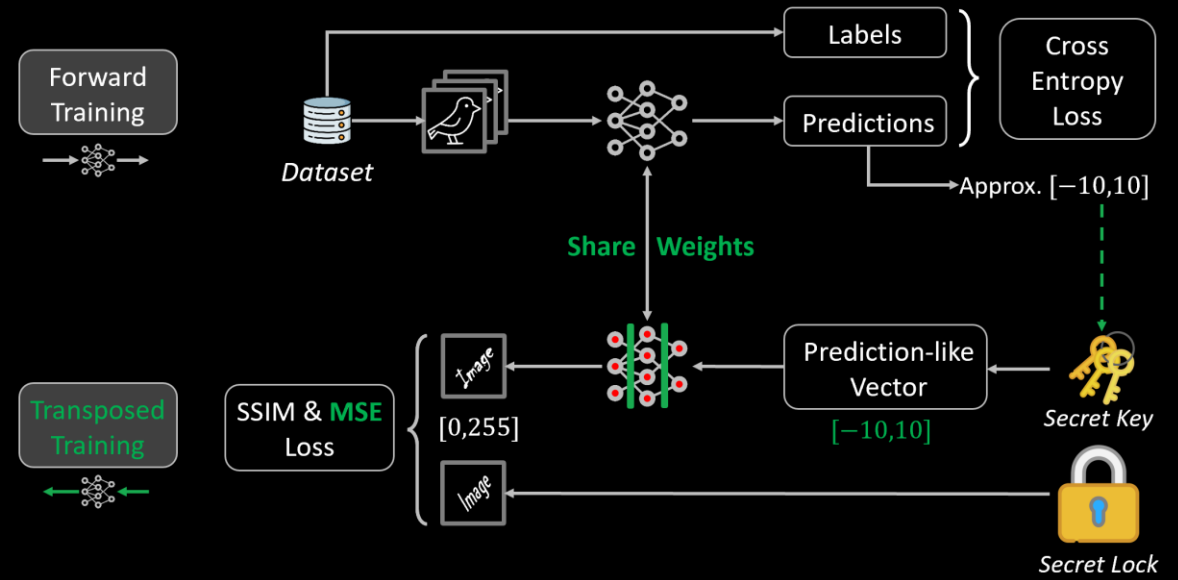


**Transposed training** to generate a **human-visible watermark**

# Thank you!!11!!1

## Any Questions?

### Parameter Entanglement



Torsten Krauß, Jasper Stang, Alexandra Dmitrienko

University of Würzburg

