



Clip-and-Verify: Linear Constraint-Driven Domain Clipping for Accelerating Neural Network Verification



Duo Zhou*, Jorge Chavez*, Hesun Chen, Grani A. Hanasusanto, Huan Zhang

*co-first authors

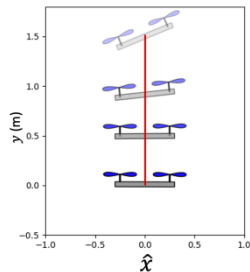
University of Illinois Urbana-Champaign



Clip-and-Verify is part of α, β -CROWN,
Winner of International Verification
of Neural Networks Competitions
(VNN-COMP 2021–2025)

The Problem: We Need Provably Safe AI

- Neural networks are used in safety-critical systems (autonomous driving, medical).
- Testing isn't enough, we need **formal proofs** of safety.
- **Challenge:** Formal verification is NP-complete and slow.
- **Goal:** Make the proof process dramatically faster.
- Existing verifiers compute a lower bound, $\underline{f}(\mathbf{x}) \leq f(\mathbf{x})$, to verify Eq. (1).



Yang et al. 2024

We must prove the quadrotor can stabilize itself from a set, \mathcal{S} , of starting points.

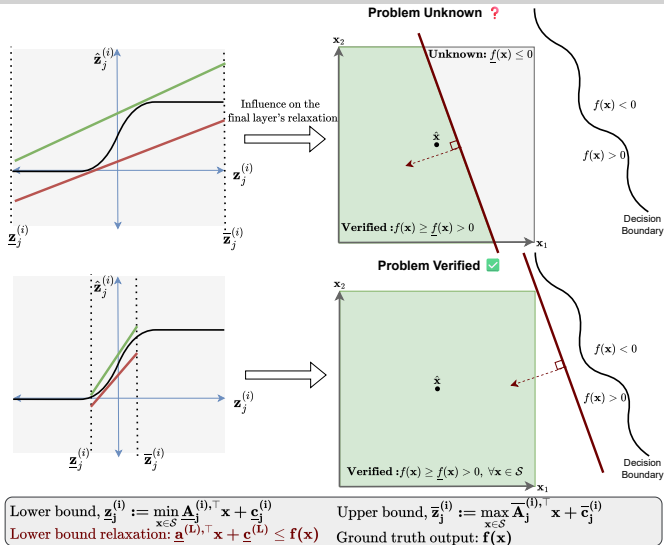
Canonical form of neural network verification

Given $\mathbf{x} \in \mathcal{S}$, is $f(\mathbf{x}) \leq 0$ satisfiable?

(1)

Contributions of Clip-and-Verify

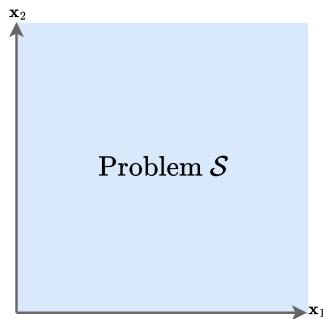
- Despite their impact on the overall relaxation, $\underline{f}(\mathbf{x})$, SOTA verifiers struggle to efficiently tighten *intermediate-layer* bounds.
- Clip-and-Verify introduces two GPU-accelerated methods to refine input domains and tighten bounds at *any* layer.



Preliminaries: How Verification Works

Branch-and-Bound (BaB): split hard regions; bound each subproblem; repeat.

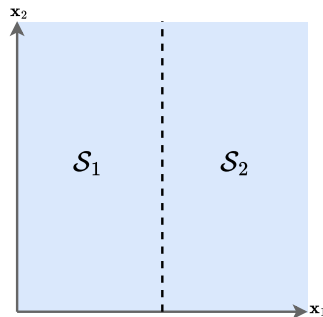
- 1 Start with one large box of inputs.
- 2 If Eq. (1) can't be verified (i.e., $\exists \mathbf{x} \in \mathcal{S}, f(\mathbf{x}) \leq 0$), **branch**: split the box.
- 3 **Bound**: compute bounds for each new box.
- 4 Recurse until the property is proved or refuted for *all* subproblems.



Preliminaries: How Verification Works

Branch-and-Bound (BaB): split hard regions; bound each subproblem; repeat.

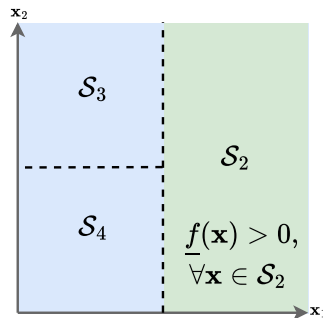
- 1 Start with one large box of inputs.
- 2 If Eq. (1) can't be verified (i.e., $\exists \mathbf{x} \in \mathcal{S}, f(\mathbf{x}) \leq 0$), **branch**: split the box.
- 3 **Bound**: compute bounds for each new box.
- 4 Recurse until the property is proved or refuted for *all* subproblems.



Preliminaries: How Verification Works

Branch-and-Bound (BaB): split hard regions; bound each subproblem; repeat.

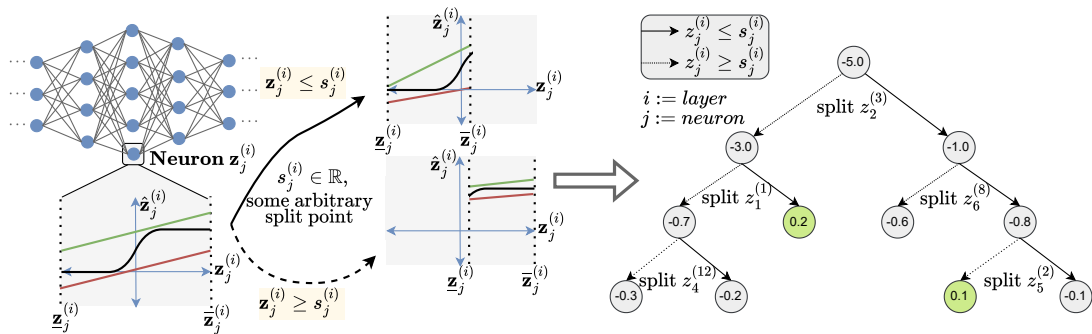
- 1 Start with one large box of inputs.
- 2 If Eq. (1) can't be verified (i.e., $\exists \mathbf{x} \in \mathcal{S}, \underline{f}(\mathbf{x}) \leq 0$), **branch**: split the box.
- 3 **Bound**: compute bounds for each new box.
- 4 Recurse until the property is proved or refuted for *all* subproblems.



Preliminaries: How Verification Works

Branch-and-Bound (BaB): Beyond the input space.

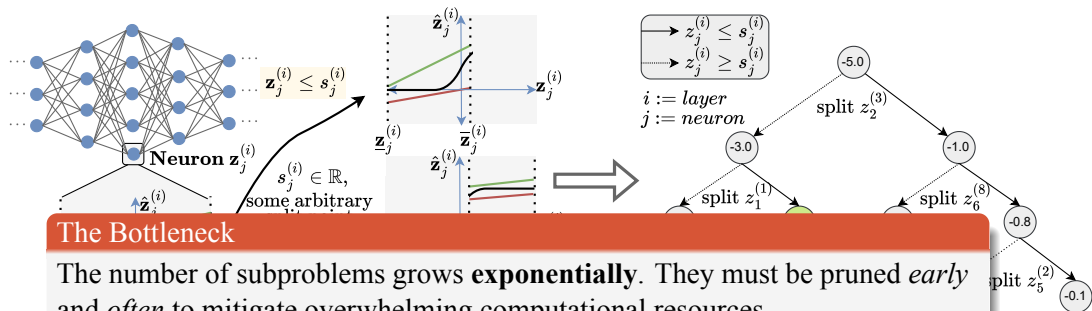
- Pre-activation neurons may be split for general activations, (e.g. $s_j^{(i)} = 0$ for ReLU) to branch off subproblems.
- Clip-and-Verify generalizes to input splits or splits on any activation function.



Preliminaries: How Verification Works

Branch-and-Bound (BaB): Beyond the input space.

- Pre-activation neurons may be split for general activations, (e.g. $s_j^{(i)} = 0$ for ReLU) to branch off subproblems.
- Clip-and-Verify generalizes to input splits or splits on any activation function.

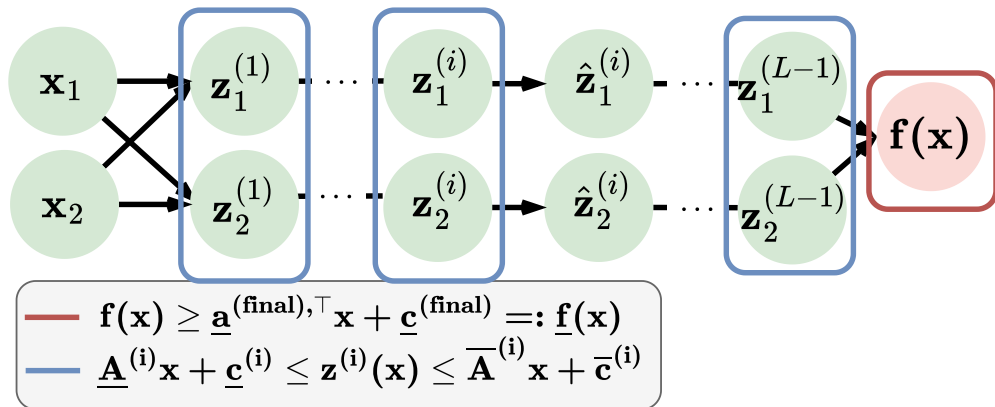


The Bottleneck

The number of subproblems grows **exponentially**. They must be pruned *early* and *often* to mitigate overwhelming computational resources.

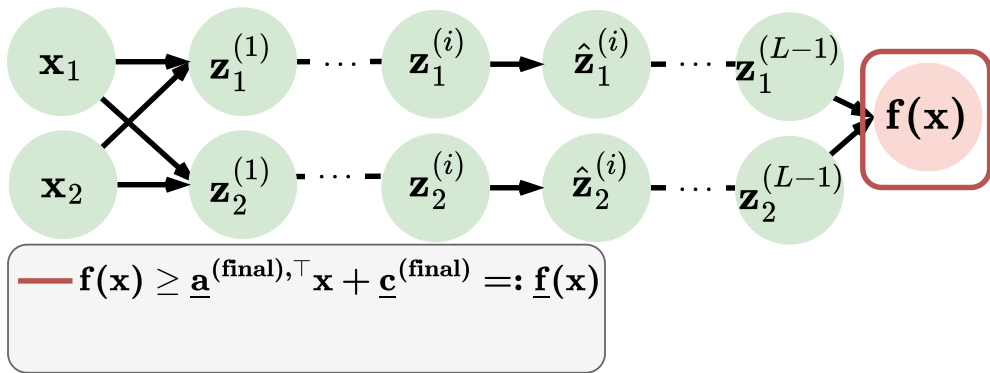
Our Core Insight: Don't Waste Information

- Linear bound propagation provides linear relaxations at every layer.



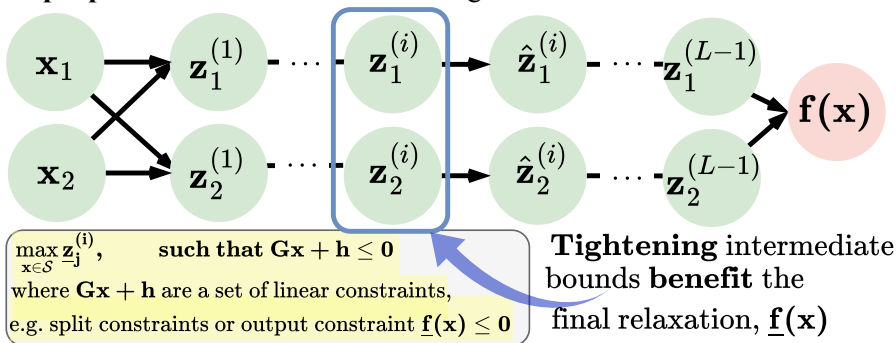
Our Core Insight: Don't Waste Information

- Linear bound propagation provides linear relaxations at every layer.
- Only the final layer's relaxation is used to determine if a subproblem has been verified.



Our Core Insight: Don't Waste Information

- Linear bound propagation provides linear relaxations at every layer.
- Only the final layer's relaxation is used to determine if a subproblem has been verified.
- Even the final layer's relaxation is **discarded** if the subproblem cannot be verified before branching further.
- Why not **re-purpose** them as constraints to strengthen the relaxation?



What are the constraints? Where do they come from?

Setting. Input box $X = \{x \mid \hat{x} - \epsilon \leq x \leq \hat{x} + \epsilon\}$.

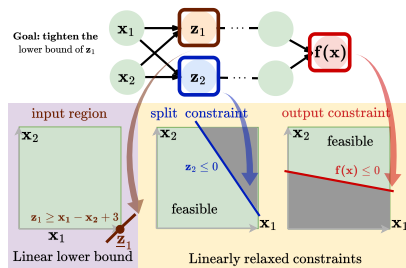
Source A - Final-layer (output) constraints.

$$f(x) \geq a^{(final),\top} x + c^{(final)} > 0 \Rightarrow a^{(final),\top} x + c^{(final)} \leq 0.$$

Source B - Activation-split constraints.

$$z_j^{(i)} \geq s \Rightarrow \bar{A}_j^{(i)\top} x + \bar{c}_j^{(i)} \geq s,$$

$$z_j^{(i)} \leq s \Rightarrow \underline{A}_j^{(i)\top} x + \underline{c}_j^{(i)} \leq s.$$



Constraints arise from BaB splits / properties.

These affine inequalities are *free byproducts*:

- tighten per-neuron bounds (Complete Clipping, 1D dual)
- shrink the input box (Relaxed Clipping, closed form)

Primal Formulation for Bound Tightening

Direct Refinement

$$L^* = \min_{\mathbf{x} \in \mathcal{S}} \{ \mathbf{a}^\top \mathbf{x} + c : \mathbf{g}^\top \mathbf{x} + h \leq 0 \} \quad (2)$$

Per-Dimension Refinement

$$\begin{aligned} \underline{x}_i &= \min_{\mathbf{x} \in \mathcal{S}} \{ x_i \mid \mathbf{A}\mathbf{x} + \mathbf{c} \leq 0 \}, \\ \bar{x}_i &= \max_{\mathbf{x} \in \mathcal{S}} \{ x_i \mid \mathbf{A}\mathbf{x} + \mathbf{c} \leq 0 \}. \end{aligned} \quad (3)$$

- It is overwhelming to refine the lower and upper bounds of *every* neuron using Eq. (2).
 - A vision transformer with $\sim 76\text{k}$ parameters amounts to solving Eq. (2) 152k times!
- Likewise, for an n -dimensional input, Eq. (3) must be solved $2n$ times.
- This is for a **single** subproblem in branch-and-bound. This is completely **intractable** for LP solvers.
- This motivates the need for an **effective** and **scalable** manner to produce such bounds.

Our Algorithms: Two Efficient Ways to Clip

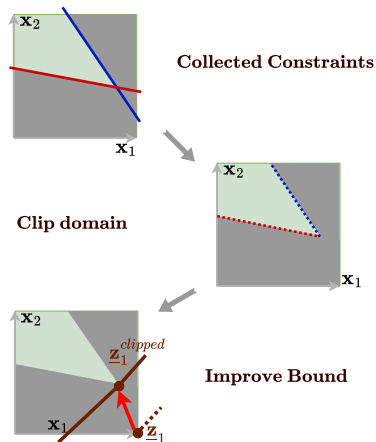
Complete Clipping:

- Tightens bounds for *specific* critical neurons.
- Solve a **1D Lagrangian dual** (not an n -D LP).
- Exact, GPU-friendly implementation.

Theorem 3.1

$$L^* = \max_{\beta \in \mathbb{R}_+} \underbrace{(a + \beta g)^\top \hat{x} - \sum_{j=1}^n |(a + \beta g)_j| \epsilon_j + c + \beta h}_{\text{Dual Objective } D(\beta)}$$

Dual $\beta \in \mathbb{R}_+$, affine relaxation, $a \in \mathbb{R}^n$, $c \in \mathbb{R}$, and linear constraint, $g^\top x + h \leq 0$, $g \in \mathbb{R}^n$, $h \in \mathbb{R}$.



Our Algorithms: Two Efficient Ways to Clip

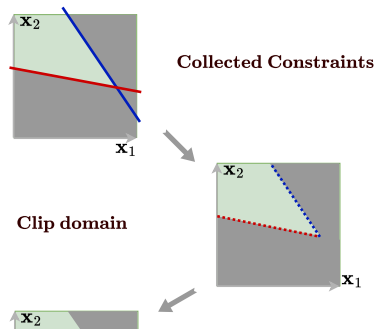
Complete Clipping:

- Tightens bounds for *specific* critical neurons.
- Solve a **1D Lagrangian dual** (not an n -D LP).
- Exact, GPU-friendly implementation.

Theorem 3.1

$$L^* = \max_{\beta \in \mathbb{R}_+} \underbrace{(a + \beta g)^\top \hat{x} - \sum_{j=1}^n |(a + \beta g)_j| \epsilon_j + c + \beta h}_{\text{Dual Objective } D(\beta)}$$

Dual $\beta \in \mathbb{R}_+$, affine relaxation, $a \in \mathbb{R}^n$, $c \in \mathbb{R}$, and linear constraint, $g^\top x + h \leq 0$, $g \in \mathbb{R}^n$, $h \in \mathbb{R}$.



Compared with LP dual solver

Specialized Complete Clipping is $\sim 880\times$ **faster** than Gurobi, a standard LP solver (comparable accuracy).

Our Algorithms: Two Efficient Ways to Clip

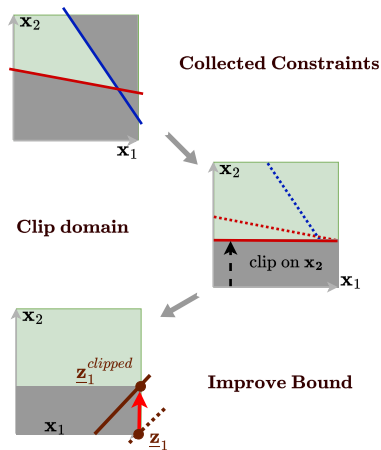
Relaxed Clipping:

- Tightens bounds for *all* neurons at once by refining the input box.
- **Closed-form** $O(n)$ update (no solver).

Theorem 3.2

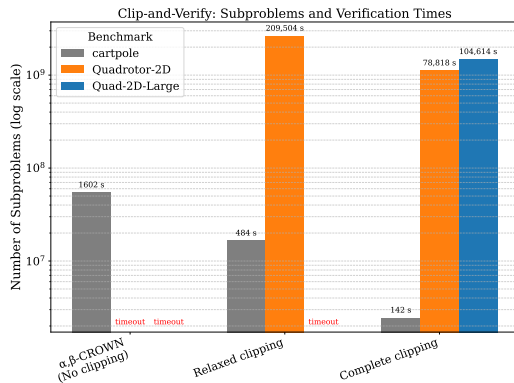
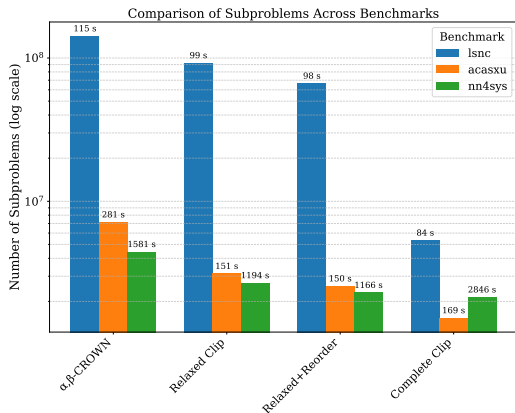
$$\begin{cases} \bar{\mathbf{x}}_i^{(\text{new})} = \min \left\{ \mathbf{x}_i^{(\text{clip})}, \bar{\mathbf{x}}_i \right\} & \text{if } a_i > 0 \\ \underline{\mathbf{x}}_i^{(\text{new})} = \max \left\{ \mathbf{x}_i^{(\text{clip})}, \underline{\mathbf{x}}_i \right\} & \text{if } a_i < 0 \\ \text{no change} & \text{otherwise} \end{cases}$$

With constraint, $\mathbf{a}^\top \mathbf{x} + c \leq 0$, $\mathbf{a} \in \mathbb{R}^n$, $c \in \mathbb{R}$, and $\mathbf{x}_i^{(\text{clip})} = (-\sum_{j \neq i} \{a_j \hat{\mathbf{x}}_j - |a_j| \epsilon_j\} - c) / a_i$.



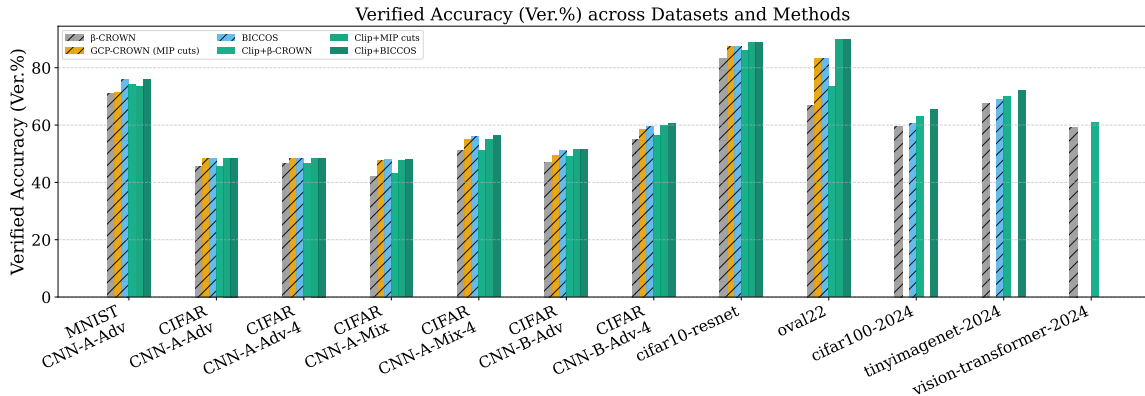
Applications: Does It Work? Yes.

- **Massive speedup:** Up to **96%** fewer BaB subproblems on `lsnc`.
- **Hard problems:** Baselines time out on control tasks (e.g., `Quadrotor-2D`); Clip-and-Verify finishes.



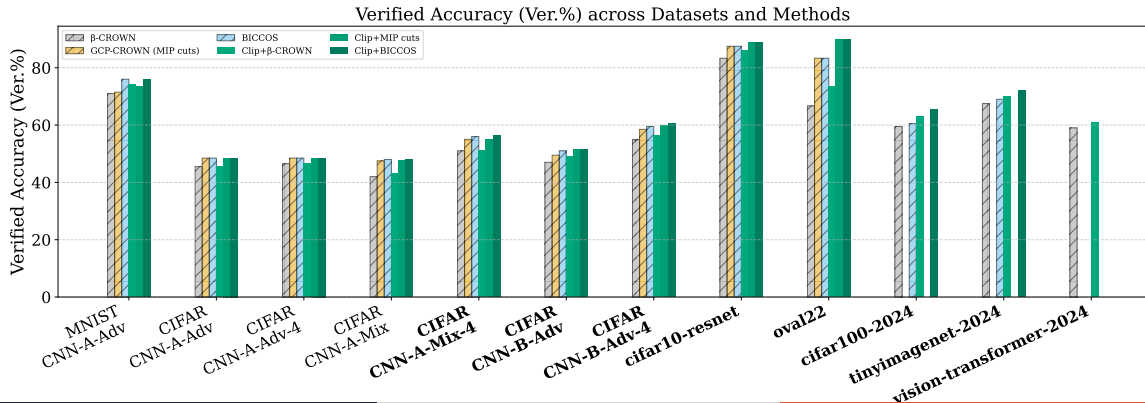
Applications: Does It Work? Yes.

- **Support all net structures:** We support any activation functions including MatMul, ReLU, Sigmoid, and Softmax.
- **SOTA:** Integrated into a VNN-COMP 2025 winning stack.



Applications: Does It Work? Yes.

- **Support all net structures:** We support any activation functions including MatMul, ReLU, Sigmoid, and Softmax.
- **SOTA:** Integrated into a VNN-COMP 2025 winning stack.



Conclusion

- **Insight:** Re-purpose **free linear constraints** that others discard.
- **Method:** Two GPU-friendly algorithms:
 - **Complete Clipping:** precise, $\sim 880\times$ faster dual solver.
 - **Relaxed Clipping:** fast, global box shrink.
- **Impact:** More properties, faster; enables previously intractable systems.
- **Code:** github.com/Verified-Intelligence/Clip_and_Verify

Thank you!



Clip-and-Verify is integrated into the state-of-the-art verifier α, β -CROWN, the overall winner of the International Verification of Neural Networks Competition (VNN-COMP) 2021–2025.