



Thesis and post-doc days, DILS
June, 17th, 2019

VERIFICATION AND VALIDATION OF DEEP LEARNING ALGORITHMS

Girard-Satabin Julien (10-2018 – 10-2021)

Chihani Zakaria (CEA)

Charpiat Guillaume, Schoenauer Marc (INRIA TAU)

PRESENTATION

Thesis and post-doc days, DILS, 2019



WHY DOING A PHD ON THIS TOPIC ?

Classical programs

```

let list_fold_right (fun x l => (x,Node proto.inputs)) ns ()
  |> fold nodes outputs ns x
let list_fold_right (fun x l => (x,Node proto.outputs)) ns ()
  |> fold outputs list_fold_right (fun x l => (x,Node proto.inputs)) ns ()
let list_tensors ns =
  let L_name x = match x, tensor_proto.name with
  | Some n -> n
  | None -> "NO NAME" in
  let list_of_n x, tensor_proto.dim in
  let L_dim x = match x, tensor_proto.raw.dim with
  | Some rd -> rd
  | None -> "NO DIM" in
  let ns = [] in
  list_fold_left (fun ns x () => match (L_name x, L_dim x, L_dim x) with
  |> ns, (L_name x, L_dim x) :: ns) ns ()
let parse_model from_file fp =
  let ch = open_in_bin fp in
  let buf = Pervasifast_list_from_channel ch in
  parse_model_proto buf
let main_graph (input proto, ti) =
  let nodes = g.nodes
  and inputs = g.inputs
  and outputs = g.outputs
  and nodes = g.nodes in
  let L_node = fold value_info.names inputs
  and n_nodes = fold value_info.names outputs
  and c_nodes = single_value_list "C_NODE" (list_length nodes) ()
  in
  let n_ops = fold nodes n_nodes
  and l_ops = ns.op (list_length l_nodes)
  and n_ops = ns.op (list_length n_nodes)
  in
  
```

Explicit control flow

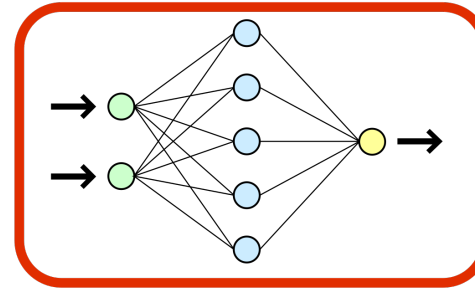
Explicit specifications

Abstractions and well known concepts

Needs to be robust

Already here

Deep learning



Generated control flow

Implicit specifications

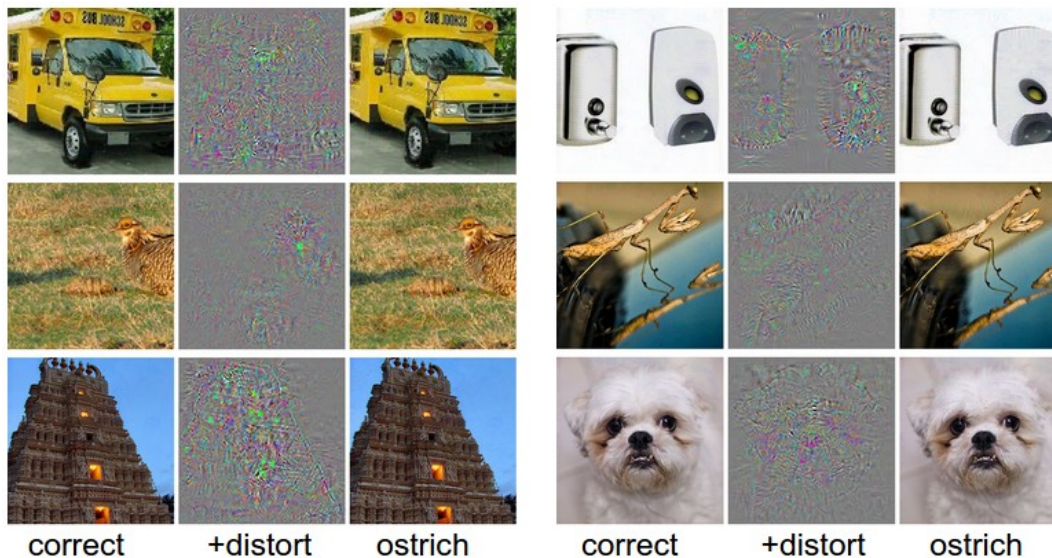
Very few abstractions and reusability

Needs to be robust

My thesis' goal



Adversarial examples



- Transferable between programs (Papernot et al.)
- No systematic way to spot them
- Easily synthesized
- Work in the physical world on image, video, sound

Vulnerabilities

- Widely focused on adversarial examples (Carlini et al., Madry et al., Goodfellow et al.,)
- Research on information leakage (Papernot et al., Tramèr et al., Abadi et al.,)
- Very few formal properties

Formally prove robustness

- Exact methods (Katz et al. for SMT-based solving, Tjeng et al. for MILP-based solutions)
- Overapproximation methods (Vechev's team with abstract interpretation, Wong et al., Tsui-wei et al., for piecewise linear overapproximations)

Vulnerabilities

- Widely focused on adversarial examples (Carlini et al., Madry et al., Goodfellow et al.,)
- Research on information leakage (Papernot et al., Tramèr et al., Abadi et al.,)
- Very few formal properties

Formally prove robustness

- Exact methods (Katz et al. for SMT-based solving, Tjeng et al. for MILP-based solutions)
- Overapproximation methods (Vechev's team with abstract interpretation, Wong et al., Tsui-wei et al., for piecewise linear overapproximations)

More can be done

- Robustness study of a new deep neural network basic block
- Bibliography
- Scientific committee (ForMaL spring school)
- Reproduction of state of the art results on adversarial robustness

- ONNX2SMT : from a neural network to a SMT formula
- Constraint-programming approaches
- Formalize privacy properties, especially differential privacy on dataset/algorithm couples

- Extend and enhance existing tools using CP
- Formalize privacy properties to check

- Course on formal methods, INRIA TAU
- Seminar on exact verification of neural networks, CEA (soon)
- Scientific committee, tutorial and presentation, ForMaL (spring school funded by DigiCOSME) : <https://formal-paris-saclay.fr/>
- *A security study of ODE Nets, Girard, Charpiat, Chihani, Schoenauer*

Commissariat à l'énergie atomique et aux énergies alternatives
Institut List | CEA SACLAY NANO-INNOV | BAT. 861 – PC142
91191 Gif-sur-Yvette Cedex - FRANCE
www-list.cea.fr

Établissement public à caractère industriel et commercial | RCS Paris B 775 685 019