

Măsurarea performanței instrumentelor specializate de verificare a rețelelor neuronale

Chira Andreea-Marina, andreea.chira00@e-uvt.ro

Hasna Diana-Mihaela, diana.hasna00@e-uvt.ro

Hasan Majd, majd.hasan93@e-uvt.ro

Chera Denis-Marian, denis.chera00@e-uvt.ro

Vanciu Catalin-Marian, catalin.vanciu00@e-uvt.ro

December 22, 2023

Abstract

Este bine cunoscut faptul că numărul de vehicule rutiere a crescut enorm datorită tehnologiilor dezvoltate din industria auto și, mai precis, disponibilității tarifelor mici. Astfel, semnele de circulație contribuie atât la siguranța șoferilor cât și a pietonilor, având un rol esențial în fluidizarea și ordonarea traficului rutier.

În acest sens, lucrarea se concentrează pe tehnicile și instrumentele disponibile pentru verificarea formală, astfel încât VNN-COMP-2023 (Concursul de verificare a rețelelor neuronale) a fost ales ca principal punct de plecare al acestui studiu. **Concursul** poate fi accesat folosind link-ul: https://github.com/ChristopherBrix/vnncomp2023_benchmarks și propune un model de rețea neuronală în formatul ONNX împreună cu o specificație în formatul VNNLIB, aceste unelte stabilind dacă există o situație în care specificația se potrivește cu modelul. Astfel de descoperiri facilitează verificarea unor proprietăți diverse în diferite aplicații ale rețelelor neuronale.

Așadar, din cea mai recentă ediție, au fost selectate două instrumente care au obținut cele mai bune performanțe, respectiv **Alpha-Beta Crown & Marabou** participând la cele mai multe criterii de referință. Prin urmare, obiectivul acestei lucrări este de a determina rezultatele obținute în urma utilizării celor două tool-uri asupra Benchmark-ului **Traffic Signs Recognition**, considerat set de date.

1 Introducere

În prezent, unul dintre cele mai utilizate sisteme de transport este reprezentat de transportul rutier, care generează un trafic semnificativ aflat într-o continuă creștere în întreaga lume, iar conform Organizației Mondiale a Sănătății (OMS) ”În fiecare an, 1,3 milioane de oameni mor și 50 de milioane sunt răniți”. Un studiu realizat de către cei

de la National Highway Transportation Safety Administration (NHTSA) din SUA în anul 2016 subliniază faptul că, procentul accidentelor rutiere datorate erorilor umane reprezintă între 94% și 96% din totalul accidentelor rutiere provocate de neglijența conducătorilor auto, ignorarea sau înțelegerea greșită a indicatoarelor rutiere, ceea ce atestă necesitatea și importanța existenței unui sistem automat și inteligent pentru a asista șoferul în sarcinile sale de conducere.

În prezent, numeroși producători auto încorporează în vehiculele lor sisteme care oferă diferite grade de autonomie și securitate. Printre aceste sisteme se numără și sistemele de recunoaștere a semnelor de circulație (TSR), care ajută la minimizarea numărului de decizii eronate și de accidente. [5]

În cazul sistemelor de tip TSR, un aspect crucial este reprezentat de rețelele neuronale care sunt folosite pentru a recunoaște pattern-uri complexe și pentru a face predicții sau decizii pe baza datelor de intrare. Rețelele neuronale sunt adesea cutii negre și este greu de asigurat că se comportă în condiții de siguranță și predictibilitate în cazul unor intrări zgomotoase sau rău intenționate (intrări malițioase). De asemenea, cercetător-ul Huan Zhang a declarat că: "Verificarea rețelelor neuronale urmărește să ofere garanții demonstrabile pentru proprietățile dorite ale rețelelor neuronale, cum ar fi siguranța și robustețea". [3] De asemenea, acestea pot fi antrenate să detecteze și să clasifice diferite tipuri de semne de circulație din imagini sau clipuri video.

Prin urmare, verificarea rețelelor neuronale urmează adesea un ciclu de verificare continuă, care implică reantrenarea rețelelor neuronale având în vedere o anumită proprietate de verificare, după cum se poate observa și în Figura 1. [4]

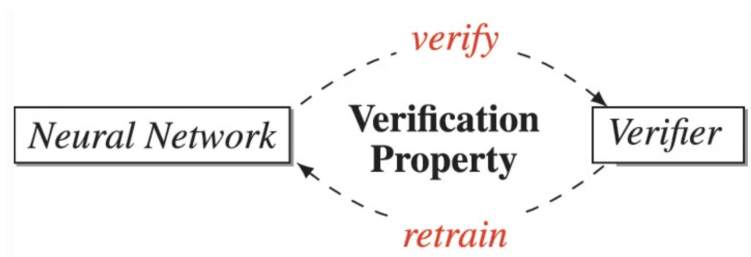


Figure 1: Ciclul de verificare a rețelelor neuronale

Astfel, în cadrul acestei lucrări, atenția este orientată spre găsirea unei modalități optime de măsurare a eficienței și corectitudinii diferitelor instrumente de verificare a rețelelor neuronale în contextul recunoașterii semnelor de circulație. Pentru a atinge acest obiectiv, vor fi detaliate, în secțiunile următoare, atât modul de funcționare al unui instrument de acest tip, cât și importanța unui benchmark și impactul pe care îl poate avea asupra performanței uneltelor selectate.

2 Benchmark

În contextul acestui experiment, s-a utilizat benchmark-ul **Traffic Sign Recognition** [1] constituit dintr-un ansamblu de fișiere ONNX (Open Neural Network Exchange) și VNNLIB, și care este folosit pentru a evalua performanța diferitelor metode și instrumente de verificare a rețelelor neuronale. Fișierele ONNX din benchmark-ul ales au la baza modelul unei rețele neuronale antrenate pe setul de date German Traffic

Sign Recognition Benchmark (GTRSB) existent pe site-ul Kaggle, care cuprinde peste 50000 de imagini ce sunt clasificate în 43 de clase, fiecare simbolizând un tip distinct de indicator rutier. Pe de altă parte, fișierele VNNLIB sunt reprezentate de problemele de verificare asociate cu acele modele, care includ diverse condiții ce pot fi legate de clasificare sau anumite proprietăți și caracteristici specifice modelelor respective.

Conținutul fișierelor VNNLIB constă în:

1. Declararea variabilelor de intrare

Fiecare valoare reprezintă o caracteristică (ex: pixel) a imaginii de intrare și sunt declarate ca numere reale.

Exemplu: (declare-const X_0 Real)

2. Constrângerile de intrare (assert-uri)

Constrângerile sunt reprezentate de condițiile impuse fiecărei variabile de intrare pentru a specifica intervalele valide pentru valorile pixelilor imaginii. Condițiile respective sunt exprimate cu ajutorul unor valori din intervalul $[0, 255]$, care în contextul imaginilor se referă la scala intensităților de lumină sau culorilor.

Exemplu:

(assert(<= X_0 255.00000000))

(assert(>= X_0 244.00000000))

3. Constrângerile de ieșire

Constrângerile de ieșire sunt definite ca o disjuncție (sau), în cadrul căreia fiecare clauză reprezintă o condiție pentru variabilele de ieșire.

Exemplu:

(assert (or
(and (>= Y_0 Y_3))
(and (>= Y_1 Y_3))))

Condițiile impun ca valorile de ieșire să fie mai mari sau egale cu valoarea de ieșire a celui de-al treilea neuron.

4. Valoarea Epsilon

Aceasta poate oscila, indicând nivelul permis de perturbație sau zgomot pentru variabilele de intrare în timpul verificării.

Exemplu: 10.00000

3 Tool

3.1 Alpha-Beta Crown

3.2 Descriere

α , β -CROWN (alpha-beta-CROWN) este un vericator de rețele neuronale care se bazează pe un mecanism eficient de propagare a limitelor liniare și pe o metodă de

”branch and bound” (Bab). Procesul poate fi accelerat în mod eficient pe GPU și se poate adapta la rețele convoluționale relativ mari (de exemplu, milioane de parametri). De asemenea, prin intermediul bibliotecii versatile auto_LiRPA acesta acceptă o gamă largă de arhitecturi de rețele neuronale, cum ar fi: CNN, ResNet și diverse funcții de activare. Totodată, tool-ul α , β -CROWN poate oferi garanții demonstrabile de robustețe împotriva atacurilor adversarilor și poate verifica, de asemenea, alte proprietăți generale ale rețelelor neuronale. [2]

Descrierea detaliată a tool-ului:

- CROWN: reprezintă un algoritm eficient de verificare bazat pe propagarea limitelor. Acesta propagă o inegalitate liniară în sens invers prin rețea și utilizează limite liniare pentru a relaxa funcțiile de activare.
- α -CROWN: este un vericator de rețele neuronale optimizat, utilizat în vericatorul Fast-and-Complete, care optimizează în comun limitele stratului intermediar și limitele stratului final în CROWN prin intermediul variabilei α . De obicei, are o putere mai mare decât LP (Linear Programming), deoarece LP nu poate restrânge în mod necostisitor limitele stratului intermediar.
- β -CROWN: introduce un parametru β optimizabil pentru a încorpora constrângerile de divizare în branch and bound (BaB) în procesul de propagare a limitelor în CROWN. Un vericator de rețele neuronale robust și scalabil este creat prin combinarea branch and bound cu o propagare eficientă și accelerată de GPU a limitelor.

3.3 Instalare

Pașii urmăți pentru configurarea tool-ului α , β -CROWN:

1. Instalare Miniconda pentru Sistemul de Operare Linux:

Pentru realizarea acestui proces s-au folosit comenzile ce pot fi regăsite în Figura 2.

```
mkdir -p ~/miniconda3
wget https://repo.anaconda.com/miniconda/Miniconda3-latest-Linux-x86_64.sh -O ~/miniconda3/miniconda.sh
bash ~/miniconda3/miniconda.sh -b -u -p ~/miniconda3
rm -rf ~/miniconda3/miniconda.sh
```

Figure 2: Instalare Miniconda

2. Clonarea proiectului α , β -CROWN

Pentru realizarea acestui proces s-a folosit comanda ce poate fi regăsită în Figura 3.

```
git clone --recursive https://github.com/Verified-Intelligence/alpha-beta-CROWN.git
```

Figure 3: Clonarea proiectului

3. Configurarea și activarea environment-ului α , β -CROWN:

Pentru parcurgerea acestei etape a fost necesară utilizarea fișierul environment.yaml aflat în folderul complete_verifier al proiectului și s-au folosit comezile ce pot fi regăsite în Figura 4.

```
# Remove the old environment, if necessary.
conda deactivate; conda env remove --name alpha-beta-crown
# install all dependents into the alpha-beta-crown environment
conda env create -f complete_verifier/environment.yaml --name alpha-beta-crown
# activate the environment
conda activate alpha-beta-crown
```

Figure 4: Environment-ul α , β -CROWN

4. Testare

După activarea environment-ului a urmat procesul de testare realizat prin intermediul comenzilor ce pot fi regăsite în Figura 5.

```
cd complete_verifier
python abcrown.py --config exp_configs/cifar_resnet_2b.yaml
```

Figure 5: Testare

În urma testării tool-ului α , β -CROWN, rezultate pot fi observate în Figura 6:

```
##### Summary #####
Initial verified acc: 36.8% (total 100 samples)
Problem instances count: 100 , total verified (safe/unsafe): 30 , total falsified (unsafe/safe): 40 , timeout: 30
mean time for all instances (total 100): 48.5506185013429, max time: 262.75381552683994
mean time for verified safe instances (total 30): 51.7745203880779, max time: 157.586641292572
mean time for verified (SAFE + UNSAFE) instances (total 70): 5.795962572097778, max time: 117.3686641292572, 8.913932542088981, 3.8509759829541, 6.6155651241335, 3.825398186479128, 3.8668753344805, 3.984576
953132328, 3.408123852585951, 3.851583565654549, 3.764037103283549, 3.8708081336493849, 3.86386122462954, 3.832833861218258, 5.423372836345431, 18.396758222778956, 3.8365688993381864, 3.815726180721084, 3.932
336248816357, 3.913855884672241, 3.865853653717041, 3.857549673583884, 3.86828558610889, 3.832977533340454, 3.786323078526123, 3.903132286241889, 3.897770456314087, 3.841413974761363, 3.864976207214355, 3.88
146045418848, 3.2385774880731
mean time for verified SAFE + TIMEOUT instances (total 60): 80.85553884364393, max time: [156.61899828918028, 122.89196721511841, 159.781444405365, 127.62200539810724, 128.1891778895922, 128.01630214718323, 137
172018313455, 136.7431731233079, 142.371948812107, 121.237265188959, 188.288815641389, 202.783855093394, 166.6924742388376, 131.82323741512842, 127.7012420881836, 138.688897168135, 122.4403731647228
1, 138.094459181585, 150.4103888188174, 117.0070524216092, 122.2382499221882, 153.5270064043396, 125.32039940055847, 134.53972148895264, 155.2627345947268, 131.62389039993286, 184.832866822052, 148.69712781
40129, 65.777489194133, 128.7656233584868]
mean time for verified UNSAFE instances (total 40): 1.314362722632682, max time: 3.55896937179554
unsafe-pgd (total 40), index: [0, 7, 8, 12, 18, 24, 27, 31, 33, 35, 40, 42, 46, 47, 48, 49, 51, 52, 53, 56, 57, 58, 59, 61, 62, 63, 65, 67, 68, 69, 76, 78, 81, 83, 85, 87, 91]
safe (total 2), index: [1, 20]
unknown (total 30), index: [2, 3, 5, 9, 11, 14, 17, 21, 23, 24, 28, 32, 36, 38, 43, 55, 64, 66, 71, 72, 77, 86, 88, 89, 93, 94, 95, 96, 97]
safe-incomplete (total 40), index: [4, 16, 38, 43, 45, 46, 49, 51, 52, 53, 54, 56, 59, 60, 64, 69, 73, 75, 76, 82, 84, 98, 99, 99]
(alpha-beta-crown) root@ujd-vicbook:~/alpha-beta-CROWN/complete_verifier [1]
```

Figure 6: Rezultatele obținute

3.4 Rularea α - β Crown pentru benchmark-ul ales

Acest proces a debutat prin clonarea proiectului VNN-COMP2023 cu scopul de a avea acces la benchmark-ul descris mai sus în cadrul acestei lucrări, și anume Traffic Sign Recognition ce are rolul de a permite evaluarea instrumentului de verificare folosit. După o căutare meticuloasă, s-a descoperit faptul că, configurațiile aferente acestei ultime ediții ale competiției se regăsesc în path-ul complete_verifier exp_configs în folderul numit vnncomp23, făcând parte din proiectul α - β Crown. Astfel, în continuare s-a identificat fișierul gtrs.yaml, care a fost configurat corespunzător pentru a stabili conexiunea necesară între tool și benchmark pentru a putea iniția, mai apoi, acțiunea de rulare propriu-zisă.

În consecință, după pregătirea configurației, pentru rulare s-a folosit comanda ce poate fi vizualizată în Figura 7, comandă cu ajutorul căreia s-a realizat testarea tool-ului adaptată astfel încat rezultatele din urma rulării să fie salvate în acest caz într-un fișier (output1.txt) pentru a putea fi vizualizate și analizate ulterior.

```
(alpha-beta-crown) root@majd-VivoBook:/# python abcrown.py --config /root/alpha-beta-CROWN/complete_verifier/exp_configs/test/gtrsb.yaml --device cpu | tee output1.txt
```

Figure 7: Rulare α , β -CROWN

Output-ul obținut în urma rulării benchmark-ului prin intermediul tool-ului poate fi observat în Figura 8 de mai jos.

```
##### Summary #####
Final verified acc: 0.0% (total 45 examples)
Problem instances count: 45 , total verified (safe/unsafe): 0 , total falsified (unsafe/sat): 43 , timeout: 2
mean time for ALL instances (total 45):346.12553217380815, max time: 6425.901163101196
mean time for verified UNSAFE instances (total 43): 75.73349313957746, max time: 2705.6366963386536
unsafe-pgd (total 42), index: [0, 1, 2, 3, 4, 7, 8, 9, 10, 11, 12, 13, 14, 15, 16, 17, 18, 19, 20, 21, 22, 23, 24, 25, 26, 27, 28, 29, 30, 31, 32, 33, 34, 36, 37, 38, 39, 40, 41, 42, 43, 44]
unknown (total 2), index: [5, 6]
unsafe-pgd (timed out) (total 1), index: [35]
```

Figure 8: Output

Rezultatele furnizate indică faptul că, din cele 45 de instanțe (fiecare reprezentând un fișier vnnlib) niciuna nu a fost verificată ca fiind sigură, majoritatea (43) au fost clasificate ca nesigure, iar două dintre cazurile nesigure au atins timeout-ul. În contextul actual, "unsafe/sat" se referă la instanțele nesigure care sunt satisfiabile, adică există cel puțin o configurație a rețelei neuronale pentru care nu a fost posibil să se demonstreze că rețeaua răspunde într-un mod dorit sau sigur.

De asemenea, timpul mediu și maxim se calculează atât pentru toate instanțele, cât și în mod particular pentru cele nesigure verificate.

References

- [1] *Traffic Signs Recognition - Benchmark*, Publicat: 05/07/2023
https://github.com/ChristopherBrix/vnncomp2023_results/tree/main/alpha_beta_crown/2023_traffic_signs_recognition
- [2] α , β -CROWN, Publicat: 10/10/2023
<https://github.com/Verified-Intelligence/alpha-beta-CROWN/blob/main/README.md>
- [3] NEURAL NETWORKS, Publicat: 15/09/2021
https://www.cs.cmu.edu/news/2021/zhang_vnn_comp2021
- [4] NETWORK ROBUSTNESS AS A VERIFICATION PROPERTY: A PRINCIPLED CASE STUDY, Publicat: 2022
https://link.springer.com/chapter/10.1007/978-3-031-13185-1_11
- [5] TRAFFIC SIGN RECOGNITION AND ANALYSIS FOR INTELLIGENT VEHICLES, Publicat: 01/03/2003
<https://www.sciencedirect.com/science/article/abs/pii/S0262885602001567>
 Link GitHub: <https://github.com/MajdHHH/Team-6-SC>