



FairQuant: Certifying and Quantifying Fairness of Deep Neural Networks

ICSE 2025: Wednesday, April 30, 2025

Brian Hyeongseok Kim*, Jingbo Wang⁺, and Chao Wang*









*University of Southern California, *Purdue University



Motivation





Public Policy





Banking

Criminal Justice





ICSE 2025: Wednesday, April 30, 2025



Training Data

Learning Algorithm

Trained Model

 $\operatorname{Does} \varphi \operatorname{hold} ?$



Motivation



Existing verifier for individual fairness [1]

- SMT solver \rightarrow limited scalability
- Qualitative \rightarrow done after finding one counterexample

Existing verifiers for local robustness [2, 3]

- Abstract interpretation \rightarrow scalable, sound
- Cannot verify individual fairness

Our Contribution: FairQuant

- Verify individual fairness of a neural network via abstract interpretation
- Quantitatively measure the degree of individual fairness across the input domain

[1] Biswas and Rajan. 2023. Fairify: Fairness Verification of Neural Networks. In 45th International Conference on Software Engineering.
[2] Wang et al. 2018. Formal security analysis of neural networks using symbolic intervals. In 27th USENIX Security Symposium.
[3] Wang et al. 2018. Efficient formal safety analysis of neural networks. In 32nd International Conference on Neural Information Processing Systems.



ICSE 2025: Wednesday, April 30, 2025





Individual Fairness



Given: classifier f, protected attribute x_i , input domain X

f satisfies <u>individual fairness</u>

with regards to x_j for the domain X

if and only if f(x) = f(x') for all $x, x' \in X$ where x and x' differ only in x_j .



(Global) Individual Fairness for the entire input domain



ICSE 2025: Wednesday, April 30, 2025

Example





Input:

x ₁ ∈ {1, 2, 3, 4, 5}
$x_2 \in \{0, 1\}$
$\mathbf{x}_3 \in \{0,1,2,3,4,5\}$

Interview score Gender, protected attribute x_j Years of experience

Output $f(x) = \{0\}$

 $f(x) = \{0, 1\}$

Binary classification



x = [1, 1, 3]	f(x) = 1
x' = [1, 0, 3]	f(x') = 0

f(x) = f(x')? X



Example





Input:

x ₁ ∈ {1, 2, 3, 4, 5}
$x_2 \in \{0, 1\}$
$\mathbf{x}_3 \in \{0,1,2,3,4,5\}$

Interview score Gender, protected attribute \boldsymbol{x}_j Years of experience

Output $f(x) = \{0,1\}$

Binary classification

100% fair or unfair?





ICSE 2025: Wednesday, April 30, 2025

Overview



1. Our Method



2. Evaluation

Dataset DNN				Fair	ify [12]				1	FairQua	nt (new)	
Dataset	DININ	Time	Cex	#Cex	Cer%	Fal%	Und%	Time	Cex	#Cex	Cer%	Fal%	Und%
	BM-1	30m	1	11	10.00	0	90.00	4.82s	1	2820	94.23	0	5.76
	BM-2	31m	1	28	16.07	0	83.93	3.23s	1	2479	93.41	0	6.58
	BM-3	31m	1	27	19.60	0	81.40	1.21s	1	1864	95.69	0	4.30
¥	BM-4	35m	1	4	3.72	0	96.18	71.12s	1	5135	87.03	0	12.96
Bar	BM-5	23m	1	114	77.25	0	22.75	1.03s	1	1474	96.27	0	3.72
	BM-6	12m	1	155	69.41	0	30.59	0.44s	1	1426	96.44	0	3.55
	BM-7	30m	1	57	9.41	0	90.59	12.26s	1	7017	83.65	0	16.34
	BM-8	30m	1	1	0.98	0	99.02	18.99s	1	3074	90.75	0	9.24
	GC-1	32m	1	22	0	0	100	9.73s	1	31585	32.67	0	67.33
=	GC-2	33m	1	6	0	0	100	31.72s	1	32655	42.21	0	57.79
ma	GC-3	8m	1	194	2.98	0	97.02	6.77s	1	25963	58.44	0	41.55
3	GC-4	4m	1	2	99.00	0	1.00	0.29s	1	77	99.65	0	0.34
	GC-5	30m	×	0	0	0	100	1.24s	1	9	99.80	0	0.19
	AC-1	32m	1	3	0.03	0	99.97	3.23s	1	6151	90.68	0	9.31
	AC-2	31m	1	9	0.01	0	99.99	30.04s	1	13008	79.93	0	20.06
	AC-3	32m	1	20	0	0	100	37.12s	1	60494	33.29	0	66.70
	AC-4	36m	×	0	0	0	100	8m	1	61324	24.79	0	75.20
	AC-5	33m	×	0	0	0	100	4m	1	71012	19.12	0	80.87
el.	AC-6	33m	1	4	0.01	0	99.99	10.20s	1	31593	58.82	0	41.17
ΡV	AC-7	30m	×	0	0.01	0	99.99	4m	1	25588	31.72	0	68.27
	AC-8	30m	1	39	0.03	0	99.97	11.18s	1	26179	66.50	0	33.49
	AC-9	30m	1	126	0.64	0	99.36	3.50s	1	5470	91.13	0	8.86
	AC-10	32m	1	8	0.03	0	99.97	5.01s	1	9033	87.65	0	12.34
	AC-11	30m	×	0	0	0	100	36.44s	1	24516	58.01	0	41.98
	AC-12	30m	×	0	0.02	0	99.98	0.91s	1	8824	70.82	0	29.17
	compas-1	17m	1	2	80.00	0.32	19.68	0.01s	1	17	97.27	2.72	0
	compas-2	31m	×	0	0	0	100	0.01s	1	15	97.59	2.40	0
as	compas-3	30m	×	0	0	0	100	0.30s	1	12	98.07	1.92	0
due	compas-4	30m	×	0	0	0	100	0.01s	1	14	97.75	2.24	0
ŭ	compas-5	T/O	×	0	0	0	100	5.24s	1	11	98.23	1.76	0
	compas-6	M/O	×	0	0	0	100	9.19s	1	12	98.07	1.92	0
	compas-7	M/O	×	0	0	0	100	101 250	1	15	07 50	2.40	0







Our Method

ICSE 2025: Wednesday, April 30, 2025



Our Method



USC Viterbi School of Engineering

ICSE 2025: Wednesday, April 30, 2025



Input: $I(x_1) = I'(x_1) = [1, 5]$ $I(x_2) = 0$ $I'(x_2) = 1$ $I(x_3) = I'(x_3) = [0, 5]$





ICSE 2025: Wednesday, April 30, 2025





Subroutine 1: Symbolic Forward Analysis Image: Constraint of the symbolic forward analysis Image: Constraint of the symbol forward analysis Input: $I(x_1) = I'(x_1) = [1, 5]$ $I(x_2) = 0$ $I'(x_2) = 1$ $I(x_3) = I'(x_3) = [0, 5]$ $i_1 \in [x_1, x_1]$ $i_1 = [1, 5]$ $I(x_2) = 0$ $I'(x_2) = 1$ $I(x_3) = I'(x_3) = [0, 5]$



 h_1

 h_2

0.5

0.7

1.2

0.4

 $i_2 \in [0,0]$ $(i_2$

 $i_3 \in [x_3, x_3]$ (i_3

0.2

-1.0

0

0.2

-1.0

 h_2

0

0.5

0.7

0.4

 $i_2 \in [1,1]$ $(i_2$

 $i_3 \in [x_3, x_3]$ (i_3













Input: $I(x_1) = I'(x_1) = [1, 5]$ $I(x_2) = 0$ $I'(x_2) = 1$ $I(x_3) = I'(x_3) = [0, 5]$





ICSE 2025: Wednesday, April 30, 2025









Input: $I(x_1) = I'(x_1) = [1, 5]$ $I(x_2) = 0$ $I'(x_2) = 1$ $I(x_3) = I'(x_3) = [0, 5]$





ICSE 2025: Wednesday, April 30, 2025















Input: $I(x_1) = I'(x_1) = [1, 5]$ $I(x_2) = 0$ $I'(x_2) = 1$ $I(x_3) = I'(x_3) = [0, 5]$



[3] Wang et al. 2018. Efficient formal safety analysis of neural networks. In 32nd International Conference on Neural Information Processing Systems.





Input: $I(x_1) = I'(x_1) = [1, 5]$ $I(x_2) = 0$ $I'(x_2) = 1$ $I(x_3) = I'(x_3) = [0, 5]$



[3] Wang et al. 2018. Efficient formal safety analysis of neural networks. In 32nd International Conference on Neural Information Processing Systems.











Input: $I(x_1) = I'(x_1) = [1, 5]$ $I(x_2) = 0$ $I'(x_2) = 1$ $I(x_3) = I'(x_3) = [0, 5]$





ICSE 2025: Wednesday, April 30, 2025













Input: $I(x_1) = I'(x_1) = [1, 5]$ $I(x_2) = 0$ $I'(x_2) = 1$ $I(x_3) = I'(x_3) = [0, 5]$



[3] Wang et al. 2018. Efficient formal safety analysis of neural networks. In 32nd International Conference on Neural Information Processing Systems.





















Our Method







Subroutine 2: Iterative Backward Analysis



 $I(x_1) = I'(x_1) = [1, 5]$ Input:

 $I(x_2) = 0$ $I'(x_2) = 1$ $I(x_3) = I'(x_3) = [0, 5]$







partition $\{x \mid x_1 \in \{1, 2, 3\}\}$



partition $\{x \mid x_1 \in \{4, 5\}\}$

Our Method







Subroutine 3: Fairness Rate Calculation



Input:

 $I(x_1) = I'(x_1) = [4, 5]$

 $I(x_2) = 0$ $I'(x_2) = 1$

 $I(x_3) = I'(x_3) = [0, 5]$









Overview



1. Our Method



2. Evaluation

Datasat DNN				Fair	ify [12]				1	FairQua	nt (new)	
Dataset	DININ	Time	Cex	#Cex	Cer%	Fal%	Und%	Time	Cex	#Cex	Cer%	Fal%	Und
	BM-1	30m	1	11	10.00	0	90.00	4.82s	1	2820	94.23	0	5.7
	BM-2	31m	1	28	16.07	0	83.93	3.23s	1	2479	93.41	0	6.5
	BM-3	31m	1	27	19.60	0	81.40	1.21s	1	1864	95.69	0	4.3
¥	BM-4	35m	1	4	3.72	0	96.18	71.12s	1	5135	87.03	0	12.9
Bar	BM-5	23m	1	114	77.25	0	22.75	1.03s	1	1474	96.27	0	3.3
	BM-6	12m	1	155	69.41	0	30.59	0.44s	1	1426	96.44	0	3.
	BM-7	30m	1	57	9.41	0	90.59	12.26s	1	7017	83.65	0	16.
	BM-8	30m	1	1	0.98	0	99.02	18.99s	1	3074	90.75	0	9.
	GC-1	32m	1	22	0	0	100	9.73s	1	31585	32.67	0	67.
=	GC-2	33m	1	6	0	0	100	31.72s	1	32655	42.21	0	57.
ma	GC-3	8m	1	194	2.98	0	97.02	6.77s	1	25963	58.44	0	41.
ð	GC-4	4m	1	2	99.00	0	1.00	0.29s	1	77	99.65	0	0.
	GC-5	30m	×	0	0	0	100	1.24s	1	9	99.80	0	0.
	AC-1	32m	1	3	0.03	0	99.97	3.23s	1	6151	90.68	0	9.
	AC-2	31m	1	9	0.01	0	99.99	30.04s	1	13008	79.93	0	20.
	AC-3	32m	1	20	0	0	100	37.12s	1	60494	33.29	0	66.
	AC-4	36m	×	0	0	0	100	8m	1	61324	24.79	0	75.
	AC-5	33m	×	0	0	0	100	4m	1	71012	19.12	0	80.
닄	AC-6	33m	1	4	0.01	0	99.99	10.20s	1	31593	58.82	0	41.
ψv	AC-7	30m	×	0	0.01	0	99.99	4m	1	25588	31.72	0	68.
	AC-8	30m	1	39	0.03	0	99.97	11.18s	1	26179	66.50	0	33.
	AC-9	30m	1	126	0.64	0	99.36	3.50s	1	5470	91.13	0	8.
	AC-10	32m	1	8	0.03	0	99.97	5.01s	1	9033	87.65	0	12.
	AC-11	30m	×	0	0	0	100	36.44s	1	24516	58.01	0	41.
	AC-12	30m	×	0	0.02	0	99.98	0.91s	1	8824	70.82	0	29.
	compas-1	17m	1	2	80.00	0.32	19.68	0.01s	1	17	97.27	2.72	
	compas-2	31m	×	0	0	0	100	0.01s	1	15	97.59	2.40	
3S	compas-3	30m	×	0	0	0	100	0.30s	1	12	98.07	1.92	
ď	compas-4	30m	×	0	0	0	100	0.01s	1	14	97.75	2.24	
ර	compas-5	T/O	×	0	0	0	100	5.24s	1	11	98.23	1.76	
	compas-6	M/O	×	0	0	0	100	9.19s	1	12	98.07	1.92	
	compas-7	M/O	x	0	0	0	100	101 250	1	15	07 50	2.40	



ICSE 2025: Wednesday, April 30, 2025

42 University of Southern California

Evaluation

FairQuant vis-a-vis Fairify [1]

32 neural networks

- 25 existing networks from Fairify \leq 300 neurons
- 7 newly trained networks up to 10000 neurons

Fairness datasets

- 1. Bank Marketing
- 2. German Credit
- 3. Adult Census

Viterbi

School of Engineering

4. Compas Recidivism

[1] Biswas and Rajan. 2023. Fairify: Fairness Verification of Neural Networks. In 45th International Conference on Software Engineering.









1. More Accurate?

2. More Scalable?

3. More Informative?



ICSE 2025: Wednesday, April 30, 2025



Fairify

20 / 32 networks found with a counterexample



32/32 networks found with a counterexample



1. More Accurate?

Method	BM-1	GC-1	AC-1	compas-1
Fairify	11	22	3	2
FairQuant	2820	31585	6151	17

of counterexamples (\uparrow) on selected models







Is FairQuant...



1. More Accurate?





2. More Scalable?

3. More Informative?



ICSE 2025: Wednesday, April 30, 2025



Fairify



Often runs until timeout Cannot verify a single partition for 3 largest networks Always finishes before timeout Often takes <1 – 100 seconds

Method	BM-1	GC-1	AC-1	compas-1
Fairify	30 min	32 min	32 min	17 min
FairQuant	4 sec	9 sec	3 sec	0.01 sec

Verification runtime (\downarrow) on selected models, T/O = 30 min



2. More Scalable?





Is FairQuant...







1. More Accurate?

2. More Scalable?

3. More Informative?



ICSE 2025: Wednesday, April 30, 2025



Fairify

Only UNSAT partitions provide Certification

FairQuant 🗸

Quantitative results for each partition



3. More Informative?

Method	BM-1	GC-1	AC-1	compas-1
Fairify	10% /	0% /	0.03% /	80% /
	90%	100%	99.97%	19.68%
FairQuant	94.23% /	32.67% /	90.68% /	97.27% /
	5.76%	67.33%	9.31%	0%

Percentages of Certified (\uparrow) / Undecided (\downarrow) rates on selected models





Is FairQuant...







1. More Accurate?

2. More Scalable?

3. More Informative?



ICSE 2025: Wednesday, April 30, 2025

Conclusion



- We introduce **FairQuant**, a novel method for verifying global individual fairness of neural networks.
- FairQuant relies on **abstract interpretation** and **iterative refinement** to achieve soundness, scalability, and accuracy.
- We provide a **quantitative** framework to measure the degree of individual fairness achieved by neural networks.



Thank You!

Any Questions?



Brian Hyeongseok Kim brian.hs.kim@usc.edu













