

Tesi di Laurea in:

# Analisi delle prestazioni di algoritmi di Machine Learning per la classificazione del traffico di rete

Relatore

Ch.mo prof. A. Pescapè

Correlatori

Ch.mo dott. M. Brescia

Ing. G. Aceto

Candidata

Luna Di Colandrea

N40/132

Anno Accademico

2012/2013



## CONTESTO

### **Analisi e monitoraggio del traffico di rete :**

Approcci di Machine Learning alla classificazione del traffico

## CONTRIBUTO

- Studio di due innovativi algoritmi di Machine Learning
  - FMLPGA (Fast Multi Layer Perceptron with Genetic Algorithm)
  - GAME (Genetic Algorithm Modeling Experiment)
- Analisi e pre-processing delle tracce di traffico a livello pacchetto utilizzate
- Progetto e esecuzione di esperimenti sui due algoritmi
- Discussione dei risultati ottenuti



# LA CLASSIFICAZIONE DEL TRAFFICO DI RETE

Con il termine **classificazione del traffico di rete** si intende l'associazione ad una sequenza di pacchetti scambiati tra due dispositivi e due rispettive porte a livello trasporto, della presunta applicazione che li ha generati

## Importanza del tema

- sicurezza informatica
- servizi differenziati
- Studio di reti nei contesti reali
- Pianificazione e gestione di risorse di sistemi di TLC

## Principali tecniche

| Port-based  | Inspection of packet payload  | Flow-based   | ML –based  |
|---|---|--|--|
| Ispezione dei 16-bit del numero di porta del livello Trasporto  | Controllo bit a bit del carico del pacchetto alla ricerca delle firme del protocollo                                      | Controllo della sola intestazione, in cui sono contenuti i valori caratteristici | Imparano da dati empirici ad associare automaticamente ad un oggetto la classe di traffico                       |
| <b>Pro:</b><br>- Semplicità   | <b>Pro:</b><br>- precisione   | <b>Pro:</b><br>- no violazione della privacy                                     | <b>Pro:</b><br>- generalizzazione,<br>- indipendenza da modelli teorici o tecniche complesse di trattamento dati |
| <b>Contro:</b><br>- porte non standardizzate<br>- Tecniche per aggirare filtri o firewall<br>- NAT<br>- FTP | <b>Contro:</b><br>- privacy<br>- crittografia e offuscamento del protocollo<br>- nuove applicazioni<br>- grandi pacchetti | <b>Contro:</b><br>- Richiede la trasmissione di tutti i pacchetti del flusso     | <b>Contro:</b><br>- Dipendenza dalla disponibilità di tracce per l'addestramento                                 |



→ algoritmi di ottimizzazione auto-adattativi

**Applicazione Classificazione:** associazione ad un campione della corrispondente etichetta (o classe)

**Le fasi di classificazione :**

- Fase di Training
- Fase di Testing
- Fase di valutazione

**Paradigmi di apprendimento**

- ML supervisionato
- ML non supervisionato

**Esempi di metodologie di apprendimento supervisionato**

Perceptron : 1 neurone di output e N neuroni di input(x)

$$y = H(\sum_{i=1}^n w_i x_i)$$

Multi-layer Perceptron : rete composta da uno o più livelli nascosti di neuroni, totalmente connessi tra i livelli di input e output

- MSE & “back propagation”



# GLI ALGORITMI GENETICI

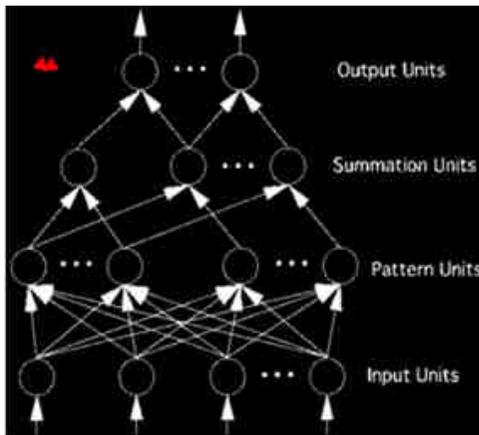
L'algoritmo genetico è un metodo di computazione ed ottimizzazione dei dati ispirato al processo evolutivo darwiniano .

## Fasi di esecuzione:

- 1) Inizializzazione della popolazione
- 2) Selezione degli output
- 3) Confronto dei risultati
- 4) Evoluzione degli individui

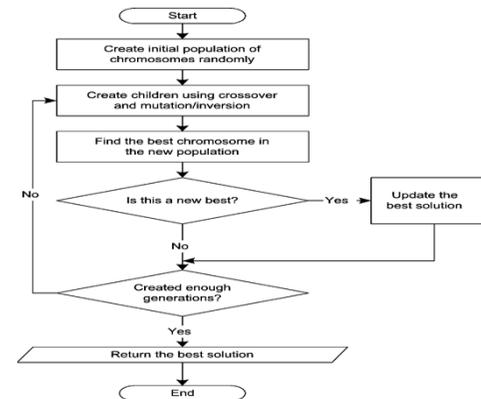
## FMLPGA

☐ l'algoritmo genetico viene impiegato come sistema di addestramento del modello MLP al posto della Back Propagation



## GAME

☐ E' un algoritmo genetico puro appositamente progettato per risolvere problemi di ottimizzazione in relazione a funzionalità di classificazione.





# TRACCE DI TRAFFICO E IL DATASET

❑ **Le tracce di flusso** : file pcap contenente 4.190.465 pacchetti , fornitoci dall'Università degli studi di Brescia<sup>1</sup> . Da questo file sono stati estratti due file : uno riguardante la Ground Truth e uno contenente i biflussi estratti .

❑ **Biflusso**: insieme di pacchetti identificati dalla stessa quintupla:

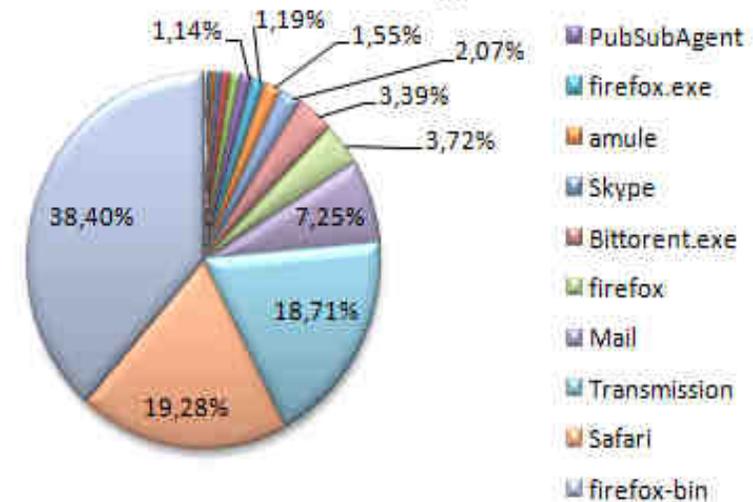
- ip sorgente
- ip destinatario
- porta sorgente
- porta destinatario
- protocollo a livello trasporto TCP/UDP

in una comunicazione **bidirezionale**.

In questo elaborato di tesi i biflussi fanno riferimento al protocollo TCP.

❑ **Costituzione del dataset**: il dataset generale è costituito da 20251 biflussi composti da 4 caratteristiche e 31 applicazioni Targets .

Percentuali Numero Biflussi Per Applicazione



❑ **Specificazione dei parametri di studio** :

- caratteristiche di input : time elapsed, byte, Up packet, Down packet
- classi Targets: applicazioni generatrici del traffico

<sup>1</sup> <http://www.ing.unibs.it/ntw/tools/traces>



## PRE-PROCESSING DEL DATASET

**Dataset utilizzati** : dalle 31 applicazioni contenute nel dataset originario sono state ricavate, per accorpamento, 6 classi principali : Browser Web, Skype, Mail, Traffico cifrato, P2P e Other. Eliminati i biflussi appartenenti alla classe Other sono stati creati 10 dataset biclasse contenenti tutti i possibili confronti delle 5 classi individuate .

**Percentuali dei dataset** : sono state usate sia le percentuali 70%(train)-30%(test) che 80%(train)-20%(test) per entrambi gli algoritmi.

**Tipi di esperimenti** : esperimenti Qualitativi (accuratezza) e Quantitativi (tempo)

## GLI ESPERIMENTI

**Scopo:** capire quali classi di traffico sono più simili tra loro e quali invece risultano più facilmente classificabili.

**Numero di prove** : per ogni esperimento eseguito sono state effettuate tre prove identiche per un totale di 189 esperimenti per GAME e 129 esperimenti per FMLPGA.

### Parametri variati:

- GAME** : grado del polinomio 5, 8, 10 e funzione di selezione degli operatori genetici Roulette, Fitting , Ranking
- **FMLPGA**: funzione di selezione degli operatori genetici Roulette(V\_1 e V\_2), Fitting, Ranking



## RISULTATI QUALITATIVI GAME

| NETWORK TRAFFIC CLASSIFICATION - TEST BASED ON GAME |    |              |         |  |       |       |       |       |       |       |       |       |       |
|---|----|--------------|---------|--|-------|-------|-------|-------|-------|-------|-------|-------|-------|
| 70% train – 30% test                                |    |              |         | <i>these columns report classification accuracy percentages (average on two classes)</i> |       |       |       |       |       |       |       |       |       |
| V   | P  | CASE         | FUNC    | Br-MI  | Br-PP | Br-Sk | Br-TC | MI-PP | MI-Sk | MI-TC | PP-Sk | PP-TC | Sk-TC |
| C<br>P<br>U   | 5  | TEST         | ROULET. | 73,75  | 63,68 | 57,82 | 68,02 | 79,43 | 75,54 | 73,62 | 62,13 | 63,43 | 73,17 |
|   |    |              | FITTING | 74,05  | 63,29 | 57,52 | 68,02 | 79,35 | 75,97 | 73,04 | 62,59 | 64,4  | 71,82 |
|   |    |              | RANK.   | 74,05  | 63,2  | 57,82 | 67,12 | 79,48 | 76,25 | 72,75 | 62,36 | 63,11 | 72,09 |
|   | 8  | TEST         | ROULET. | 81,95  | 66,76 | 57,23 | 61,71 | 85,93 | 87,84 | 80    | 69,84 | 63,75 | 72,9  |
|   |    |              | FITTING | 81,06  | 67,44 | 55,46 | 63,96 | 85,5  | 87,41 | 80,87 | 70,75 | 63,43 | 72,63 |
|   |    |              | RANK.   | 81,43  | 67,34 | 59    | 60,81 | 84,91 | 87,7  | 80,87 | 70,52 | 61,17 | 72,63 |
|   | 10 | TEST         | ROULET. | 81,58  | 66,47 | 59,29 | 64,86 | 85,25 | 86,84 | 80,58 | 68,03 | 64,08 | 73,98 |
|   |    |              | FITTING | 81,8   | 67,44 | 58,7  | 66,22 | 85,46 | 88,41 | 83,19 | 68,25 | 62,78 | 74,25 |
|   |    |              | RANK.   | 82,48  | 66,57 | 58,7  | 65,32 | 85,38 | 87,41 | 80    | 69,55 | 63,43 | 72,63 |
|   |    | NUM PATTERNS | 447     | 346  | 113   | 74    | 791   | 233   | 115   | 147   | 103   | 123   |       |



# RISULTATI QUALITATIVI GAME

| NETWORK TRAFFIC CLASSIFICATION - TEST BASED ON GAME |     |     |   |         |       |       |       |       |       |       |       |       |       |       |
|---|-----|-----|---|---------|-------|-------|-------|-------|-------|-------|-------|-------|-------|-------|
| 70% train – 30% test                                |     |     | these columns report classification accuracy percentages (average on two classes) |         |       |       |       |       |       |       |       |       |       |       |
| V   | P   | CAS | FUNC  | B-MI    | Br-PP | Br-Sk | Br-TC | MI-PP | MI-Sk | MI-TC | PP-Sk | PP-TC | Sk-TC |       |
| 5   | TES | T   | ROULET.   | 73,75   | 63,68 | 57,82 | 68,02 | 79,43 | 75,54 | 73,62 | 62,13 | 63,43 | 73,17 |       |
|   |     |     | FITTING   | 74,05   | 63,29 | 57,52 | 68,02 | 79,35 | 75,97 | 73,04 | 62,59 | 64,4  | 71,82 |       |
|   |     |     | RANK.   | 74,05   | 63,2  | 57,82 | 67,12 | 79,48 | 76,25 | 72,75 | 62,36 | 63,11 | 72,09 |       |
|   | 8   | TES | T   | ROULET. | 81,95 | 66,76 | 57,23 | 61,71 | 85,93 | 87,84 | 80    | 69,84 | 63,75 | 72,9  |
|   |     |     |   | FITTING | 81,06 | 67,44 | 55,46 | 63,96 | 85,5  | 87,41 | 80,87 | 70,75 | 63,43 | 72,63 |
|   |     |     |   | RANK.   | 81,43 | 67,34 | 59    | 60,81 | 85,91 | 87,7  | 80,6  | 70,52 | 61,7  | 72,63 |
| 10  | TES | T   | ROULET.   | 81,58   | 66,47 | 59,29 | 60,86 | 85,25 | 86,84 | 80,58 | 68,9  | 64,08 | 73,98 |       |
|   |     |     | FITTING   | 81,8    | 67,44 | 58,7  | 62,22 | 85,46 | 88,41 | 83,19 | 68,7  | 62,78 | 74,25 |       |
|   |     |     | RANK.   | 82,48   | 66,57 | 58,7  | 65,32 | 85,38 | 87,41 | 80    | 69,55 | 63,43 | 72,63 |       |
| NUM PATTERNS  |     |     |   | 447     | 346   | 113   | 74    | 791   | 233   | 115   | 147   | 103   | 123   |       |

|         | MI-PP | MI-Sk | MI-TC |
|---------|-------|-------|-------|
| ROULET. | 85,25 | 86,84 | 80,58 |
| FITTING | 85,46 | 88,41 | 83,19 |
| RANK.   | 85,38 | 87,41 | 80    |
|         | 791   | 233   | 115   |

| Matrice di Confusione caso Mail-Skype |          |          |                |     |     |
|---------------------------------------|----------|----------|----------------|-----|-----|
| TRAINING                              |          | TEST     |                |     |     |
|                                       | P(C1 Sk) | P(C2 MI) |                |     |     |
| target class 1                        | 219      | 47       | target class 1 | 103 | 18  |
| target class 2                        | 15       | 260      | target class 2 | 7   | 105 |

|                        |          |
|------------------------|----------|
| total classification   | 88,54%   |
| class 1 classification | 82,3308% |
| Class 2 classification | 94,5455% |

|                        |        |
|------------------------|--------|
| total classification   | 89,27% |
| class 1 classification | 85,12% |
| Class 2 classification | 93,75% |

Il dataset con i migliori risultati è Mail-Skype



## RISULTATI QUALITATIVI FMLPGA

| NETWORK TRAFFIC CLASSIFICATION - TEST BASED ON FMLPGA |              |         |   |       |       |       |       |       |       |       |       |       |
|---|--------------|---------|---|-------|-------|-------|-------|-------|-------|-------|-------|-------|
| 80% train -- 20% test                                 |              |         | <i>classification accuracy percentages (average on two classes)</i> |       |       |       |       |       |       |       |       |       |
| V   | CASE         | FUNC    | Br-MI   | Br-PP | Br-Sk | Br-TC | MI-PP | MI-Sk | MI-TC | PP-Sk | PP-TC | Sk-TC |
| C<br>P<br>U   | TEST         | ROULV1  | 91,61   | 78,79 | 91,89 | 79,59 | 93,74 | 94,19 | 93,51 | 78,57 | 71,43 | 91,46 |
|   |              | ROULV2  | 91,95   | 81,82 | 90,54 | 75,51 | 93,36 | 93,55 | 92,21 | 86,73 | 87,14 | 92,68 |
|   |              | FITTING | 89,6  | 71,43 | 91,89 | 73,47 | 90,89 | 90,32 | 93,51 | 79,59 | 72,86 | 92,68 |
|   |              | RANKING | 89,93   | 79,65 | 91,89 | 81,63 | 91,46 | 93,55 | 92,21 | 85,71 | 81,43 | 91,46 |
|   | NUM PATTERNS | 298     | 231   | 74    | 49    | 527   | 155   | 77    | 98    | 70    | 82    |       |
| G<br>P<br>U   | TEST         | ROULV1  | 90,27   | 79,65 | 97,3  | 81,63 | 92,6  | 92,9  | 94,81 | 93,88 | 71,43 | 91,46 |
|   |              | ROULV2  | 91,61   | 81,82 | 91,89 | 85,71 | 91,84 | 94,84 | 93,51 | 83,67 | 68,57 | 92,68 |
|   |              | FITTING | 90,27   | 80,95 | 90,54 | 79,59 | 92,03 | 94,84 | 93,51 | 92,86 | 70    | 92,68 |
|   |              | RANKING | 91,28   | 80,82 | 93,15 | 84,1  | 91,45 | 92,74 | 92,98 | 87,62 | 70,32 | 90,14 |
|   | NUM PATTERNS | 298     | 231   | 74    | 49    | 527   | 155   | 77    | 98    | 70    | 82    |       |



# RISULTATI QUALITATIVI FMLPGA

| NETWORK TRAFFIC CLASSIFICATION - TEST BASED ON FMLPGA |      |         |  |           |           |           |           |           |           |           |           |           |
|---|------|---------|--|-----------|-----------|-----------|-----------|-----------|-----------|-----------|-----------|-----------|
| 80% train -- 20% test                                 |      |         | classification accuracy percentages (average on two classes) |           |           |           |           |           |           |           |           |           |
| V   | CASE | FUNC    | Br-MI  | Br-PP     | Br-Sk     | Br-TC     | MI-PP     | MI-Sk     | MI-TC     | PP-Sk     | PP-TC     | Sk-TC     |
| C   | TEST | ROULV1  | 91,6<br>1  | 78,7<br>9 | 91,8<br>9 | 79,5<br>9 | 93,7<br>4 | 94,1<br>9 | 93,5<br>1 | 78,5<br>7 | 71,4<br>3 | 91,4<br>6 |
|   |      | ROULV2  | 91,9<br>5  | 81,8<br>2 | 90,5<br>4 | 75,5<br>1 | 93,3<br>6 | 93,5<br>5 | 92,2<br>1 | 86,7<br>3 | 87,1<br>4 | 92,6<br>8 |
|   |      | FITTING | 89,6<br>3  | 71,4<br>3 | 91,8<br>9 | 73,4<br>7 | 90,8<br>9 | 90,3<br>2 | 93,5<br>1 | 79,5<br>9 | 72,8<br>6 | 92,6<br>8 |
|   |      | RANKING | 89,9<br>3  | 79,6<br>5 | 91,8<br>9 | 81,6<br>3 | 91,4<br>6 | 93,5<br>5 | 92,2<br>1 | 85,7<br>1 | 81,4<br>3 | 91,4<br>6 |
| NUM PATTERNS  |      |         | 298  | 231       | 74        | 49        | 527       | 155       | 77        | 98        | 70        | 82        |
| G   | TEST | ROULV1  | 90,2<br>7  | 79,6<br>5 | 97,3<br>3 | 81,6<br>3 | 92,6<br>3 | 92,9<br>9 | 94,8<br>1 | 93,8<br>8 | 71,4<br>3 | 91,4<br>6 |
|   |      | ROULV2  | 91,6<br>1  | 81,8<br>2 | 91,8<br>9 | 85,5<br>4 | 91,8<br>4 | 94,8<br>4 | 93,5<br>1 | 83,6<br>7 | 68,5<br>7 | 92,6<br>8 |
|   |      | FITTING | 90,2<br>7  | 80,9<br>5 | 90,5<br>4 | 75,5<br>3 | 92,0<br>3 | 94,8<br>4 | 93,5<br>1 | 92,8<br>6 | 70,3<br>7 | 92,6<br>8 |
|   |      | RANKING | 91,2<br>8  | 80,8<br>2 | 93,1<br>5 | 84,1<br>3 | 91,4<br>3 | 92,7<br>4 | 92,9<br>8 | 87,6<br>2 | 80,3<br>2 | 90,1<br>4 |
| NUM PATTERNS  |      |         | 298  | 231       | 74        | 49        | 527       | 155       | 77        | 98        | 70        | 82        |

|         | MI-PP | MI-Sk | MI-TC | PP-Sk |
|---------|-------|-------|-------|-------|
| ROULV1  | 92,6  | 92,9  | 94,81 | 93,88 |
| ROULV2  | 91,84 | 94,84 | 93,51 | 83,67 |
| FITTING | 92,03 | 94,84 | 93,51 | 92,86 |
| RANKING | 91,45 | 92,74 | 92,98 | 87,62 |
|         | 527   | 155   | 77    | 98    |

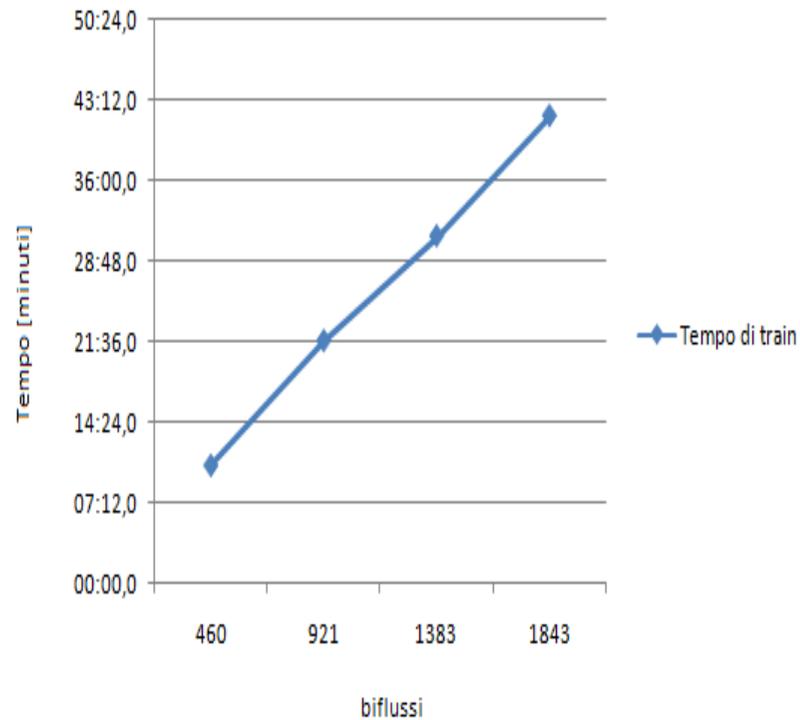
Le percentuali di classificazione risultano più omogenee rispetto a GAME



# RISULTATI QUANTITATIVI

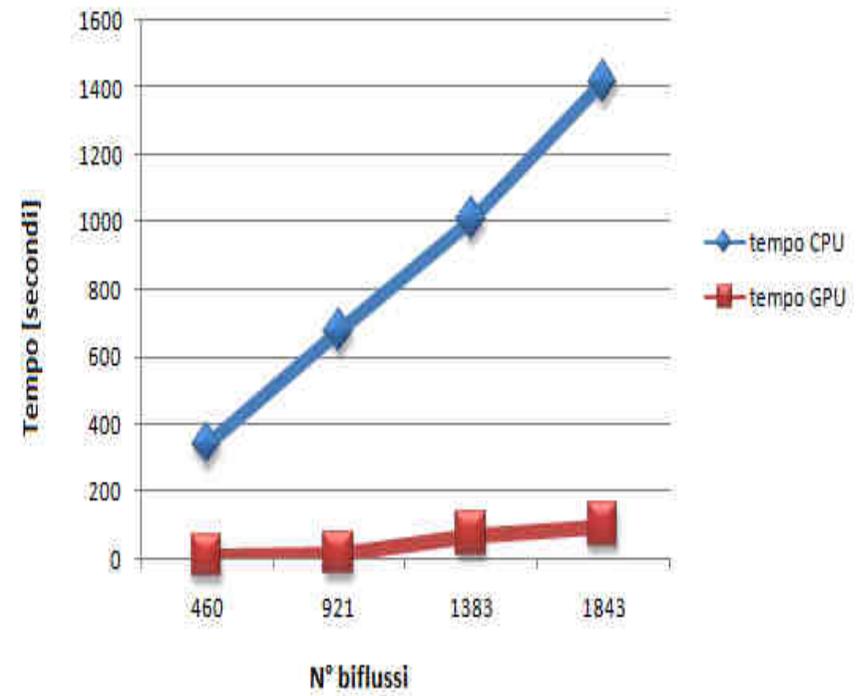
## GAME

Andamento del tempo di train al variare del numero di biflussi



## FMLPGA

Andamento del tempo di Train al variare del numero di biflussi



Speed Up : 25x



## CONCLUSIONI

In questo elaborato di tesi abbiamo:

- introdotto un pre-processing delle tracce di traffico per poter utilizzare gli algoritmi di classificazione considerati
- analizzato le prestazioni di due algoritmi di ML per la classificazione del traffico di rete

## LAVORI FUTURI

- Esecuzione di esperimenti di tipo multiclasse
- Studio delle prestazioni di GAME su architetture GPU
- Esperimenti su altre tracce di traffico contenenti anche altre classi di traffico

**Grazie per  
l'attenzione**