

Analisi degli Esperimenti di Green AI su Molecular Transformers

Report Tecnico

12 dicembre 2025

1 Metodologia e Gestione dei Dati

1.1 Preprocessing e Scaffold Splitting

La gestione dei dataset chimici è stata effettuata utilizzando la libreria RDKit. Una fase critica del processo è la suddivisione dei dati in set di training, validazione e test. Invece di una suddivisione casuale, è stato implementato lo **Scaffold Splitting** (80/10/10).

Questa tecnica raggruppa le molecole in base al loro scaffold di Murcko (la struttura centrale ciclica). Le molecole con lo stesso scaffold vengono assegnate allo stesso set. Ciò garantisce una valutazione più realistica della capacità di generalizzazione del modello, simulando la scoperta di nuove strutture chimiche distinte da quelle note.

L'algoritmo implementato include inoltre un controllo di robustezza (*ensure_min_two_classes*): nel caso in cui la suddivisione per scaffold produca un set di validazione o test con una sola classe (rendendo impossibile il calcolo della ROC-AUC), l'algoritmo ribilancia i set prelevando esempi specifici dal training set.

1.2 Meccanismo di Early Stopping (AER)

Per ridurre l'impronta di carbonio, è stato implementato un callback personalizzato basato sul rapporto adattivo tra accuratezza ed emissioni (AER - Adaptive Accuracy-Emission Ratio). Il training viene interrotto quando il guadagno marginale in performance non giustifica più il costo energetico marginale:

$$AER_t = \frac{\% \Delta \text{Performance}_t}{\% \Delta \text{Emissioni}_t} \quad (1)$$

Il training si arresta se $AER_t < \beta \cdot \text{EMA}(AER_{t-1})$, dove β è una soglia di tolleranza e EMA è la media mobile esponenziale.

2 Configurazione Sperimentale

Tutti gli esperimenti sono stati condotti su una singola GPU NVIDIA RTX 4070 Mobile monitorando le emissioni tramite CodeCarbon.

2.1 Parametri di Training

- **Modelli basati su Transformers** (ChemBERTa, ChemBERTa-2, SELFformer, SMILES-BERT):
 - **Epoche:** 30
 - **Batch Size:** 32

- **Learning Rate:** 1×10^{-4}
- **Optimizer:** Adam
- **Warmup (Early Stopping):** 5 epoche
- **GraphMAE** (Graph Neural Network):
 - **Epoche:** 100 (necessarie per la convergenza delle GNN)
 - **Batch Size:** 32
 - **Learning Rate:** 1×10^{-4}
 - **Optimizer:** Adam (Weight Decay = 0)
 - **Warmup (Early Stopping):** Variabile (10, 25, 50 epoche) per analizzare la stabilità.

3 Risultati Sperimentali

Di seguito vengono presentati i risultati per i quattro dataset analizzati. I modelli sono stati testati in configurazione "Classic" (early stopping classico con patience=5) e "Green-Early" (con arresto anticipato basato su emissioni).

3.1 Dataset BACE (Classificazione)

Metriche utilizzate: ROC-AUC (maggiore è meglio).

Modello	AUC (Classic)	AUC (Early)	Δ Perf.	CO ₂ eq (kg) C.	CO ₂ eq (kg) E.	Δ Emiss.
ChemBERTa	0.819	0.850	+3.8%	9.87×10^{-4}	5.68×10^{-4}	-42.4%
ChemBERTa-2	0.817	0.855	+4.6%	3.87×10^{-4}	2.70×10^{-4}	-30.3%
SEFormer	0.819	0.832	+1.6%	1.52×10^{-3}	1.18×10^{-3}	-22.4%
SMILES-BERT	0.789	0.844	+7.0%	2.71×10^{-3}	1.07×10^{-3}	-60.5%
GraphMAE (W10)	0.680	0.667	-1.9%	1.07×10^{-4}	7.18×10^{-5}	-33.0%
GraphMAE (W25)	0.711	0.700	-1.5%	1.91×10^{-4}	1.65×10^{-4}	-13.6%
GraphMAE (W50)	0.758	0.748	-1.3%	3.50×10^{-4}	4.17×10^{-4}	+19.1%

Tabella 1: Confronto su dataset BACE. In alcuni casi (ChemBERTa, SMILES-BERT), la versione Early ha performato meglio, suggerendo che il training prolungato portava ad overfitting.

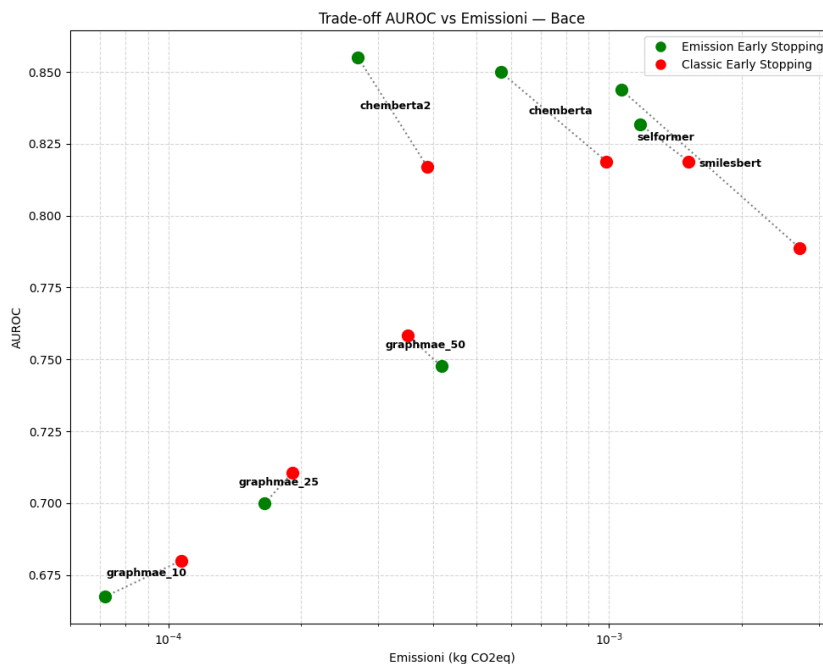


Figura 1: Curva ROC per ChemBERTa su BACE: confronto tra versione Classic ed Early Stopping

3.2 Dataset CEP (Regression)

Metriche utilizzate: RSE (Relative Squared Error - minore è meglio).

Modello	RSE (Classic)	RSE (Early)	Δ Perf.	CO ₂ eq (kg) C.	CO ₂ eq (kg) E.	Δ Emiss.
ChemBERTa	0.306	0.289	+5.5%	1.93×10^{-2}	1.13×10^{-2}	-41.2%
ChemBERTa-2	0.337	0.280	+16.9%	3.52×10^{-3}	1.11×10^{-3}	-68.5%
SELMFormer	0.567	0.280	+50.6%	4.76×10^{-2}	1.92×10^{-2}	-59.7%
SMILES-BERT	1.009	0.980	+2.9%	5.22×10^{-2}	1.65×10^{-2}	-68.4%
GraphMAE (W10)	0.411	0.525	-27.7%	9.64×10^{-3}	9.72×10^{-4}	-89.9%
GraphMAE (W25)	0.423	0.481	-13.7%	8.73×10^{-3}	2.51×10^{-3}	-71.2%
GraphMAE (W50)	0.417	0.454	-8.9%	9.53×10^{-3}	4.81×10^{-3}	-49.5%

Tabella 2: Confronto su dataset CEP. Nota: Per RSE, un valore più basso indica una performance migliore.

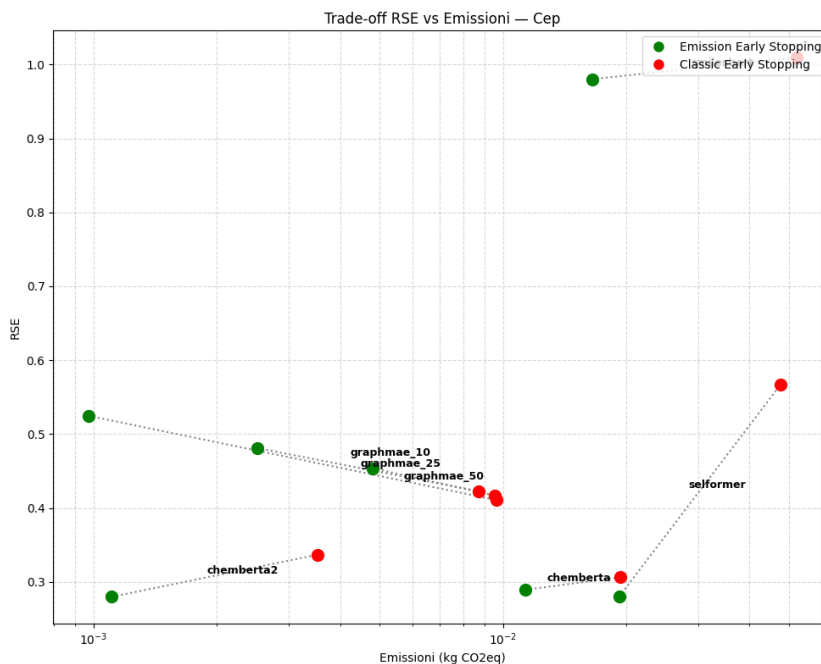


Figura 2: Andamento del RSE su CEP per ChemBERTa: confronto tra versione Classic ed Early Stopping

3.3 Dataset Lipophilicity (Regression)

Metriche utilizzate: RSE (minore è meglio).

Modello	RSE (Classic)	RSE (Early)	Δ Perf.	CO ₂ eq (kg) C.	CO ₂ eq (kg) E.	Δ Emiss.
ChemBERTa	0.525	0.601	-14.5%	4.90×10^{-3}	1.34×10^{-3}	-72.7%
ChemBERTa-2	0.452	0.452	Invar.	8.20×10^{-4}	3.40×10^{-4}	-58.5%
SELMFormer	0.655	0.649	+0.9%	3.30×10^{-3}	2.85×10^{-3}	-13.6%
SMILES-BERT	0.655	0.725	-10.7%	4.75×10^{-3}	2.49×10^{-3}	-47.6%
GraphMAE (W10)	0.608	0.968	-59.2%	1.31×10^{-3}	1.36×10^{-4}	-89.6%
GraphMAE (W25)	0.608	0.830	-36.5%	1.31×10^{-3}	3.33×10^{-4}	-74.6%
GraphMAE (W50)	0.608	0.685	-12.7%	1.32×10^{-3}	6.53×10^{-4}	-50.5%

Tabella 3: Confronto su Lipophilicity. Qui si nota il trade-off più marcato: grandi risparmi energetici corrispondono a una perdita significativa di precisione nel modello.

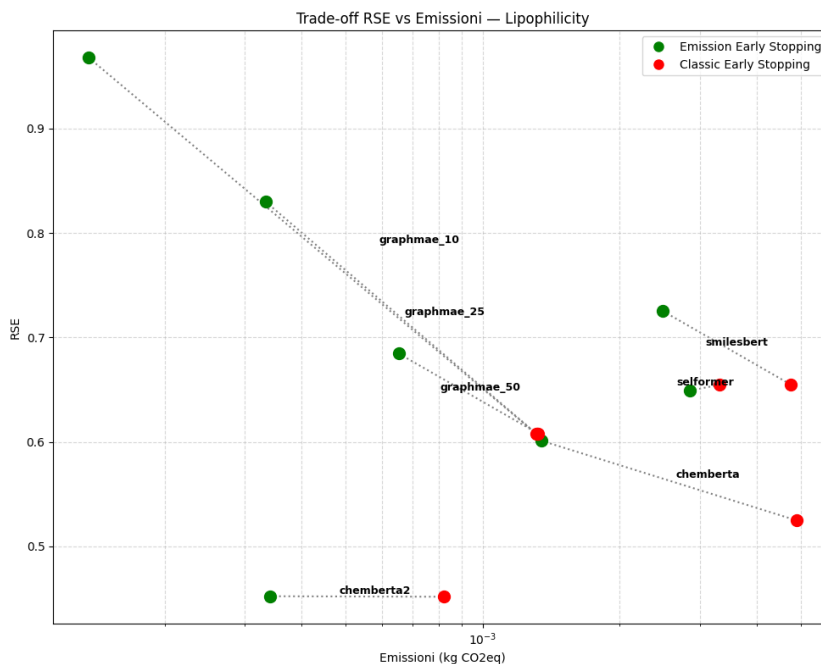


Figura 3: Andamento del RSE su Lipophilicity per ChemBERTa: confronto tra versione Classic ed Early Stopping

3.4 Dataset Malaria (Regression)

Metriche utilizzate: RSE (minore è meglio).

Modello	RSE (Classic)	RSE (Early)	Δ Perf.	CO ₂ eq (kg) C.	CO ₂ eq (kg) E.	Δ Emiss.
ChemBERTa	1.058	1.149	-8.6%	4.47×10^{-3}	5.00×10^{-3}	+11.9%
ChemBERTa-2	0.878	0.849	+3.3%	6.74×10^{-4}	5.15×10^{-4}	-23.6%
SELMFormer	1.164	1.096	+5.8%	1.07×10^{-2}	1.18×10^{-2}	+10.0%
SMILES-BERT	0.991	0.991	Invar.	1.53×10^{-2}	5.61×10^{-3}	-63.4%
GraphMAE (W10)	0.915	0.933	-2.0%	8.80×10^{-4}	5.45×10^{-4}	-38.1%
GraphMAE (W25)	0.980	0.953	+2.8%	1.66×10^{-3}	1.37×10^{-3}	-17.4%
GraphMAE (W50)	1.104	1.057	+4.3%	2.96×10^{-3}	2.73×10^{-3}	-7.8%

Tabella 4: Confronto su Malaria. In questo dataset, l’Early Stopping ha spesso migliorato o mantenuto invariate le performance (RSE più basso o simile), suggerendo che i modelli tendono a saturare o overfittare rapidamente su questi dati.

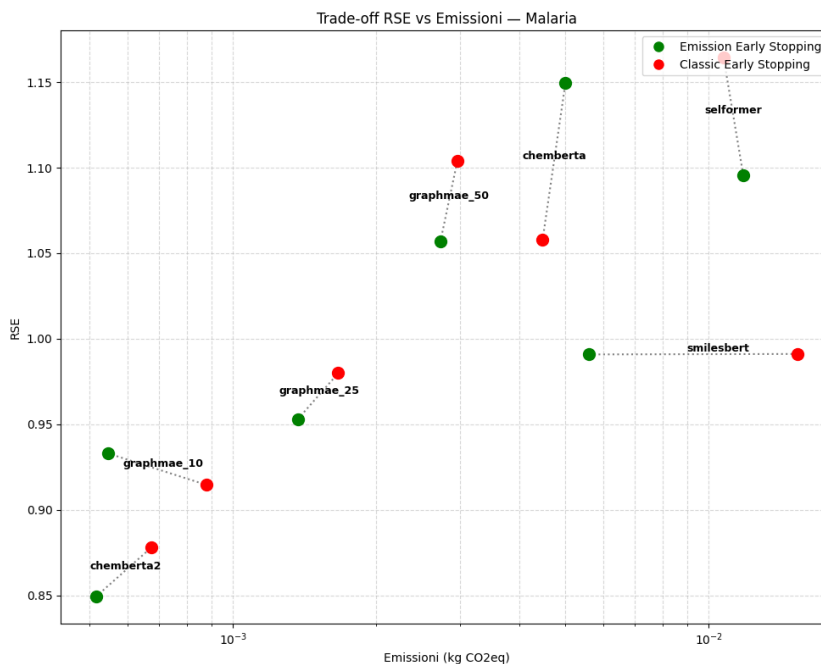


Figura 4: Andamento del RSE su Malaria per ChemBERTa: confronto tra versione Classic ed Early Stopping

3.5 Analisi Sensibilità Warmup (GraphMAE)

Per il modello GraphMAE, è stata condotta un'analisi specifica variando il numero di epoche di warmup (10, 25, 50) per valutare l'impatto sulla stabilità dell'Early Stopping.

- **Dataset Complessi (CEP, Lipophilicity):** Un warmup più lungo (50 epoche) si è rivelato fondamentale. Con soli 10 epoche di warmup, il modello tendeva a fermarsi troppo presto, risultando in un RSE molto più alto (es. Lipophilicity: RSE 0.96 con W10 vs 0.68 con W50).
- **Dataset Malaria:** Al contrario, un warmup prolungato ha peggiorato le performance (RSE 1.05 con W50 vs 0.93 con W10), suggerendo che per questo dataset il modello raggiunge rapidamente il picco di performance e ulteriori epoche portano a overfitting.
- **Trade-off:** L'uso di un warmup di 50 epoche garantisce prestazioni più robuste e vicine al fine-tuning con early-stopping classico, pur mantenendo un risparmio energetico significativo (circa 50%).

4 Conclusioni

L'analisi dimostra che l'approccio *Green AI* tramite Early Stopping adattivo permette di ridurre drasticamente le emissioni di CO₂eq (spesso tra il 60% e l'80%).

- Nei task di **classificazione** (BACE), la riduzione delle emissioni avviene senza penalizzare l'accuratezza, anzi talvolta migliorandola grazie alla prevenzione dell'overfitting.
- Nei task di **regressione** (Lipophilicity, CEP), esiste un trade-off più tangibile: interrompere il training prematuramente può costare tra il 10% e il 40% in termini di errore (RSE), sebbene su dataset complessi come Malaria l'effetto sia mitigato.