



THE UNIVERSITY
OF QUEENSLAND
AUSTRALIA

CREATE CHANGE

Reduce, Reuse, Recycle: Green Information Retrieval Research

Harry Scells, Shengyao Zhuang, Guido Zucccon

h.scells@uq.edu.au

The University of Queensland, Australia



NLP

ML

[1] Strubell, E. et al. 2019. Energy and Policy Considerations for Deep Learning in NLP. *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*



NLP

ML

What about IR research?

[1] Strubell, E. et al. 2019. Energy and Policy Considerations for Deep Learning in NLP. *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*

But what are emissions?

- **Energy:** *amount of work done*
 - Measured in **joules**

But what are emissions?

- **Energy:** *amount of work done*
 - Measured in **joules**
- **Power:** *energy per unit time*
 - Measured in **watts**; 1 watt = 1 joule/second
 - kWh: energy consumed at a rate of 1 kilowatt for 1 hour

But what are emissions?

- **Energy:** *amount of work done*
 - Measured in **joules**
- **Power:** *energy per unit time*
 - Measured in **watts**; 1 watt = 1 joule/second
 - kWh: energy consumed at a rate of 1 kilowatt for 1 hour
- **Emissions:** *by-products created by producing power*
 - Measured in kgCO₂e; kilograms of carbon dioxide equivalent



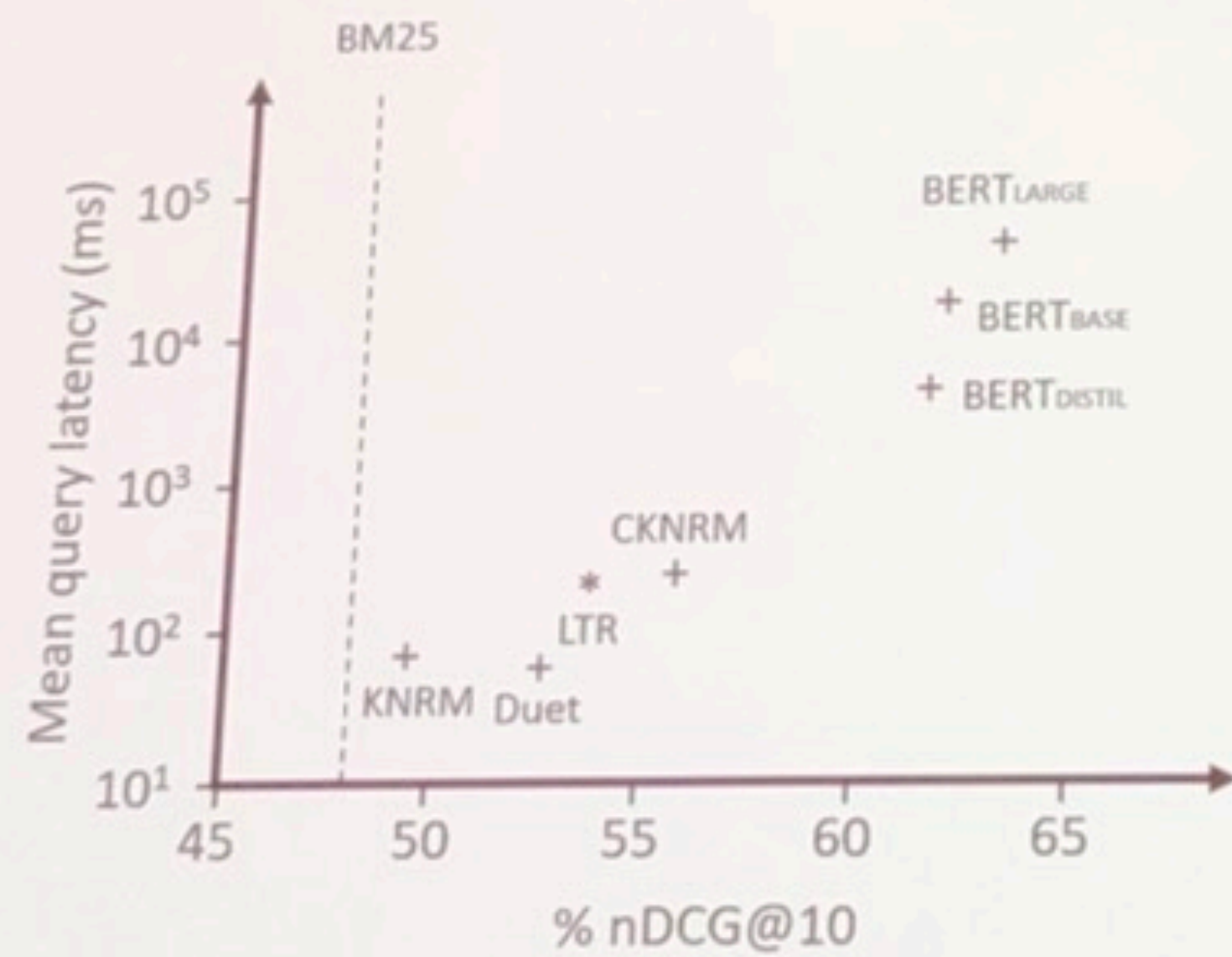
NLP

ML

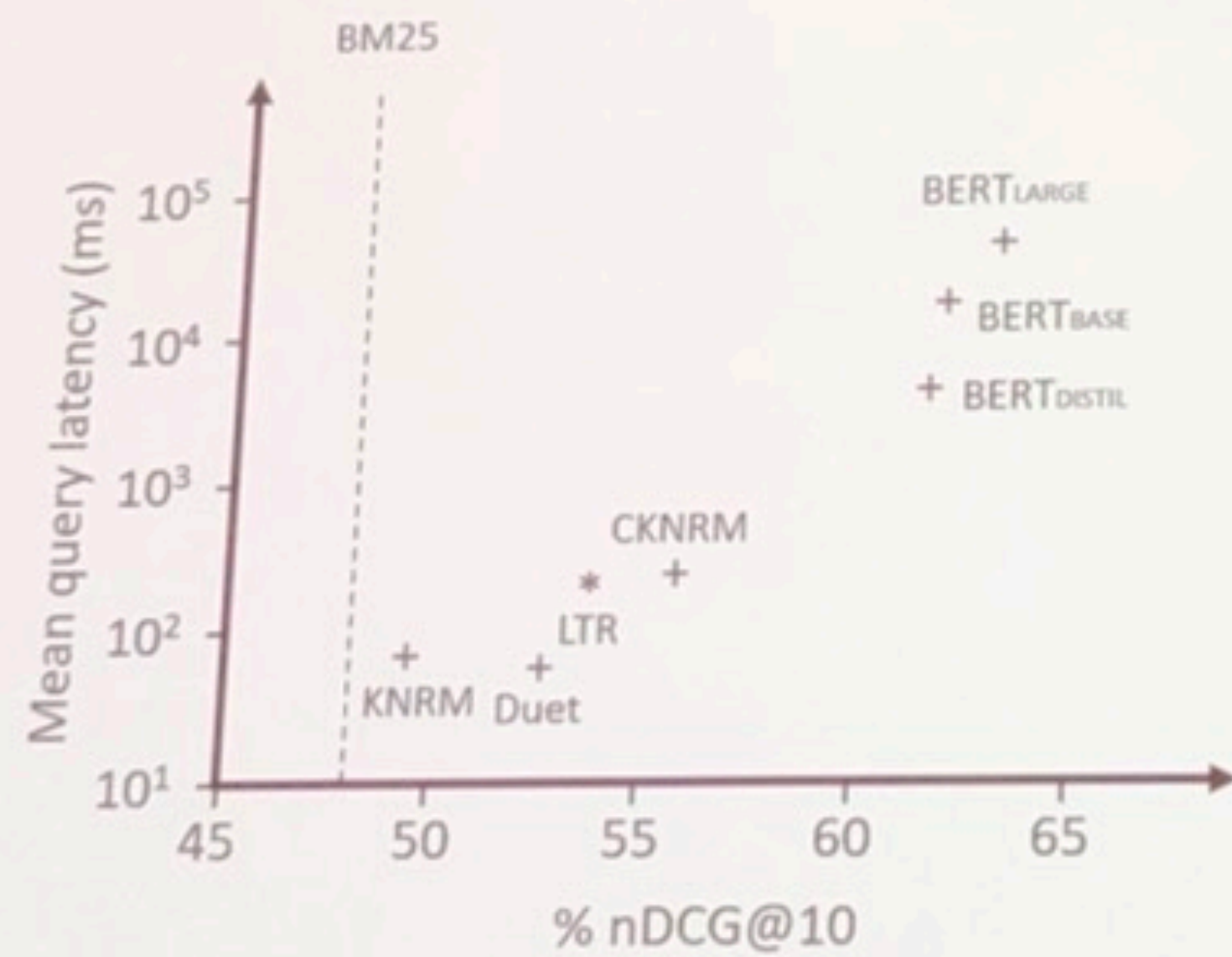
What about IR research?

Isn't this just retrieval efficiency?

[1] Strubell, E. et al. 2019. Energy and Policy Considerations for Deep Learning in NLP. *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*

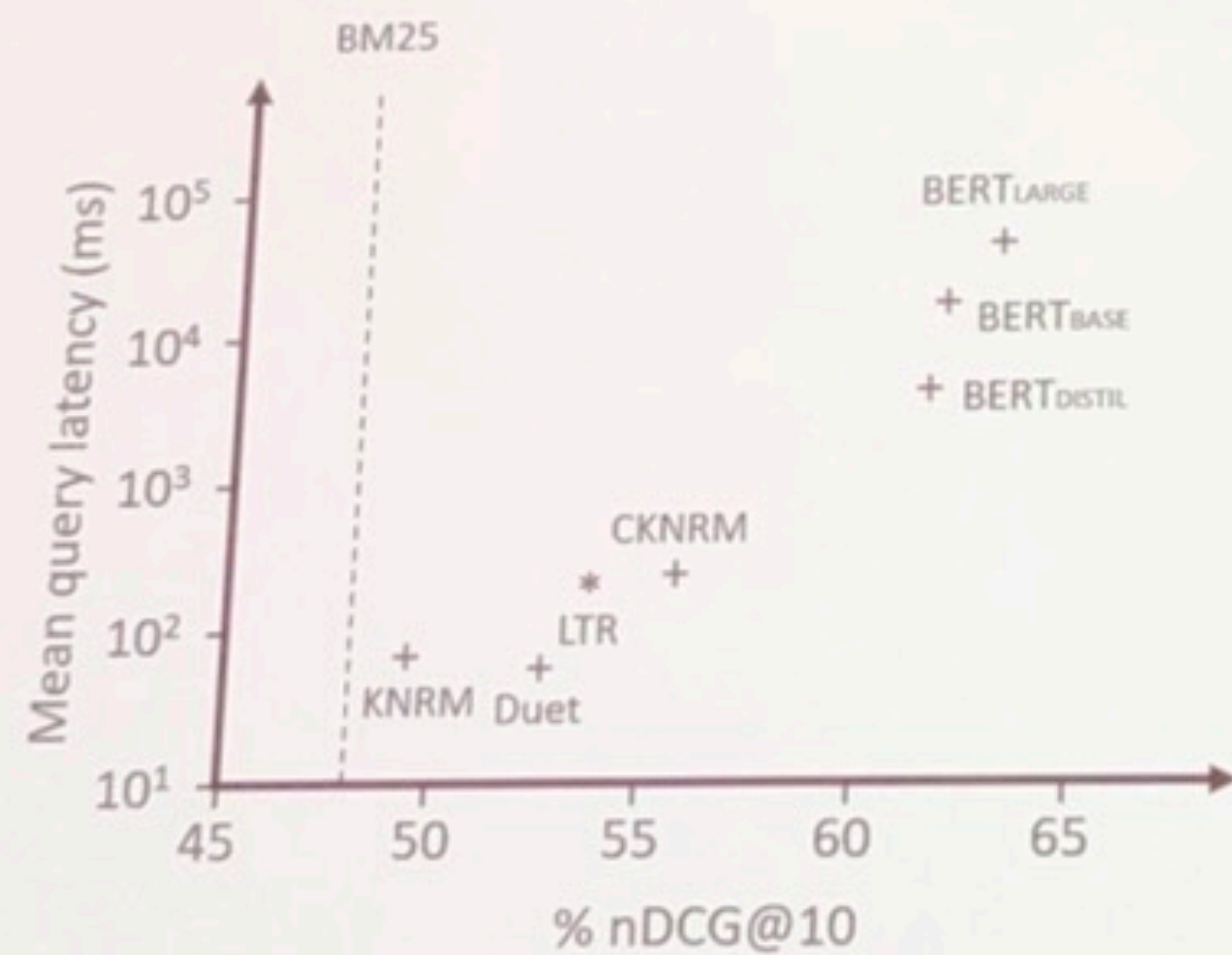


TREC Deep Learning Track 2019 – Document Ranking

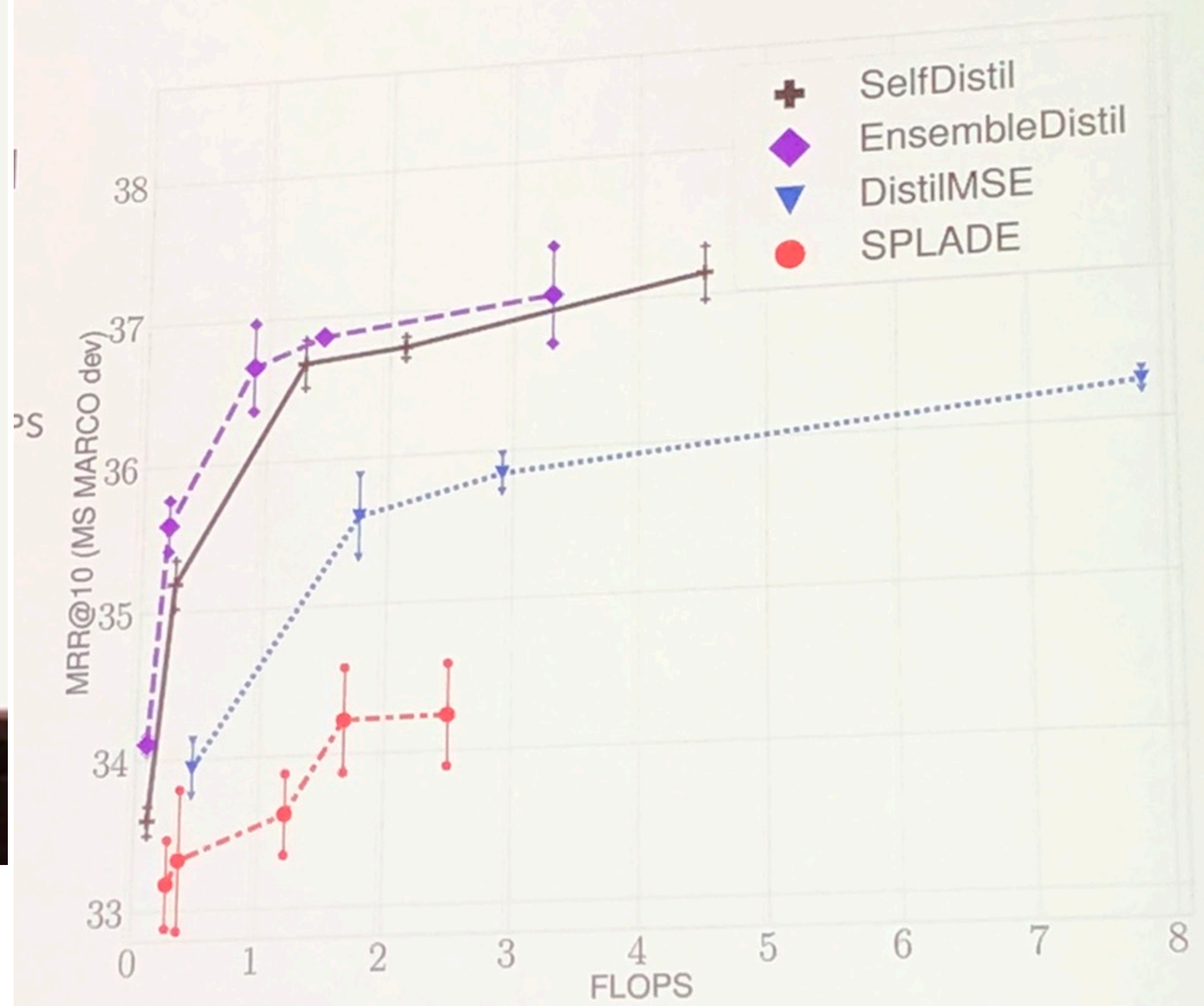
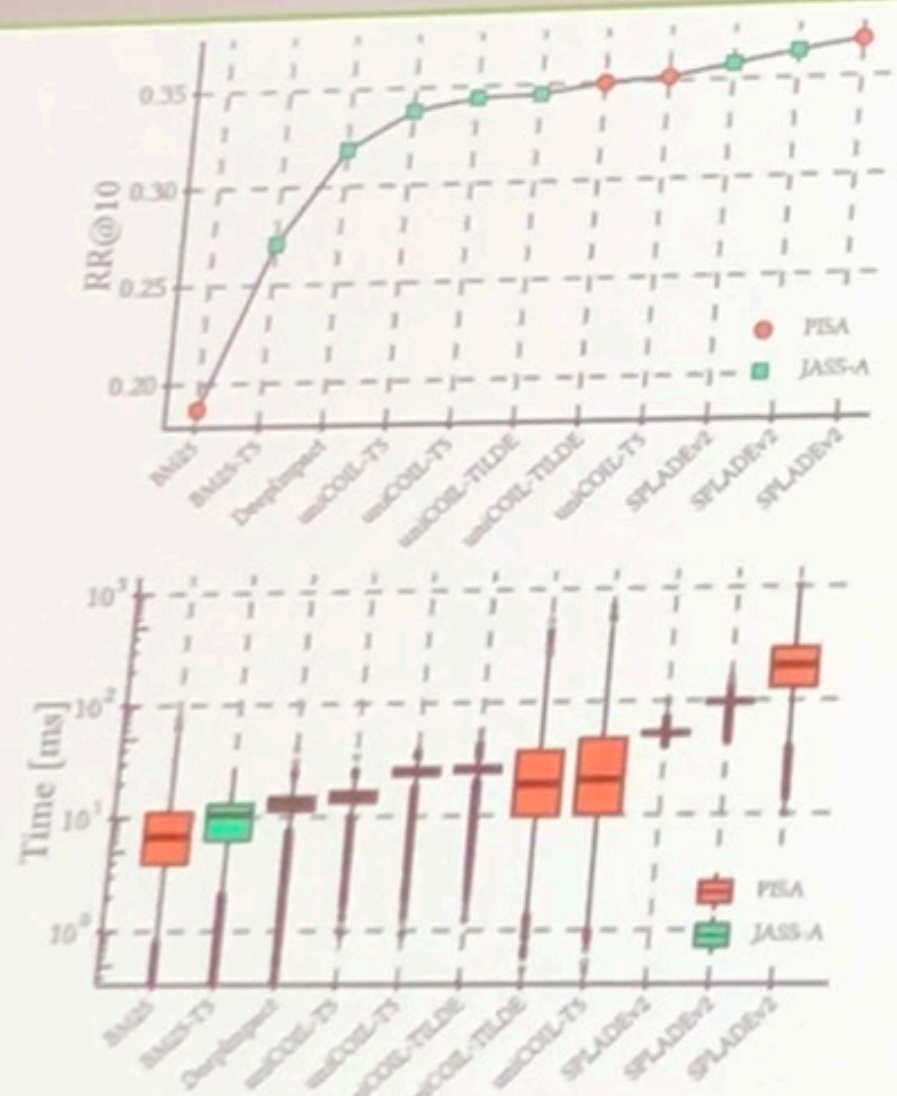


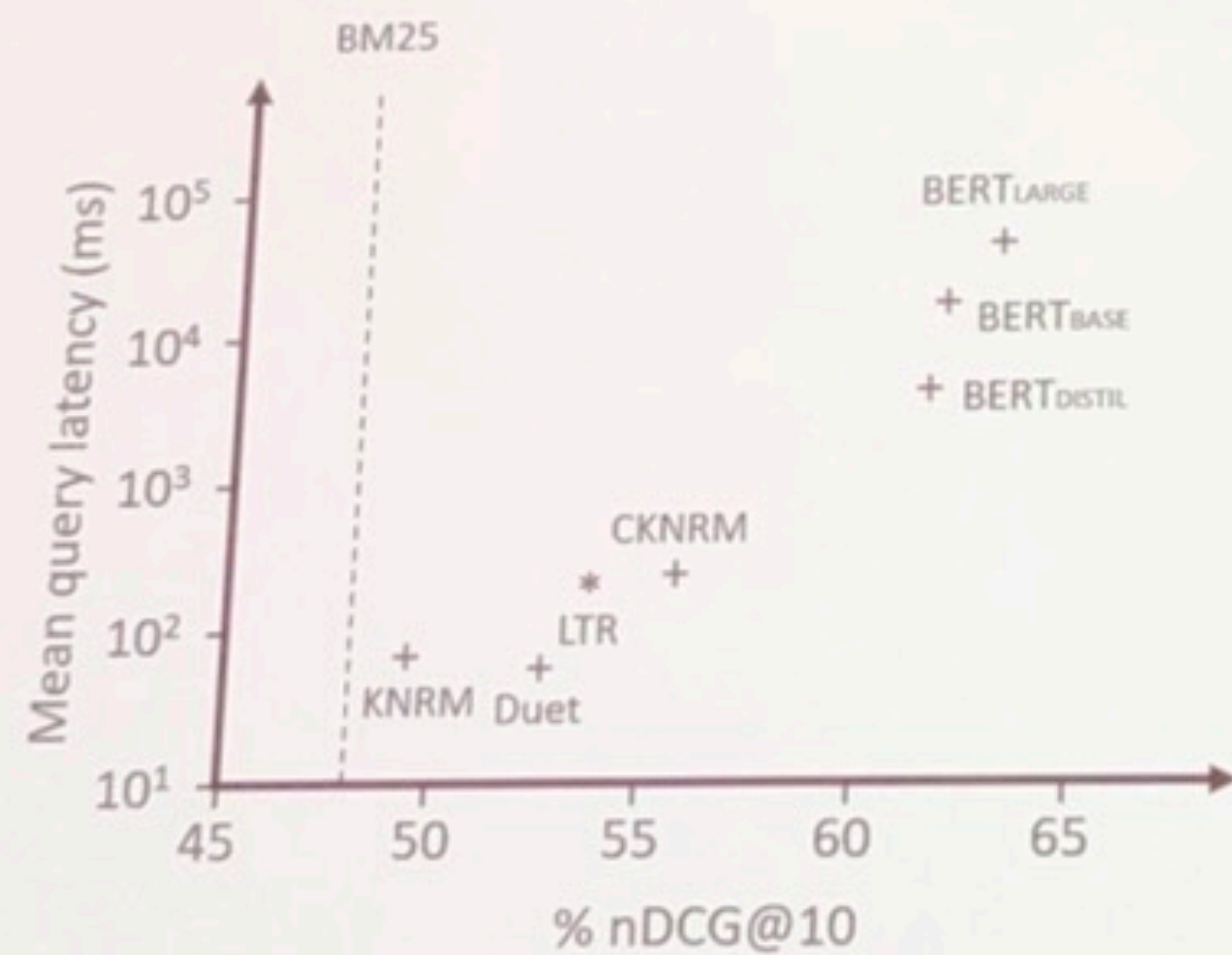
TREC Deep Learning Track 2019 – Document Ranking



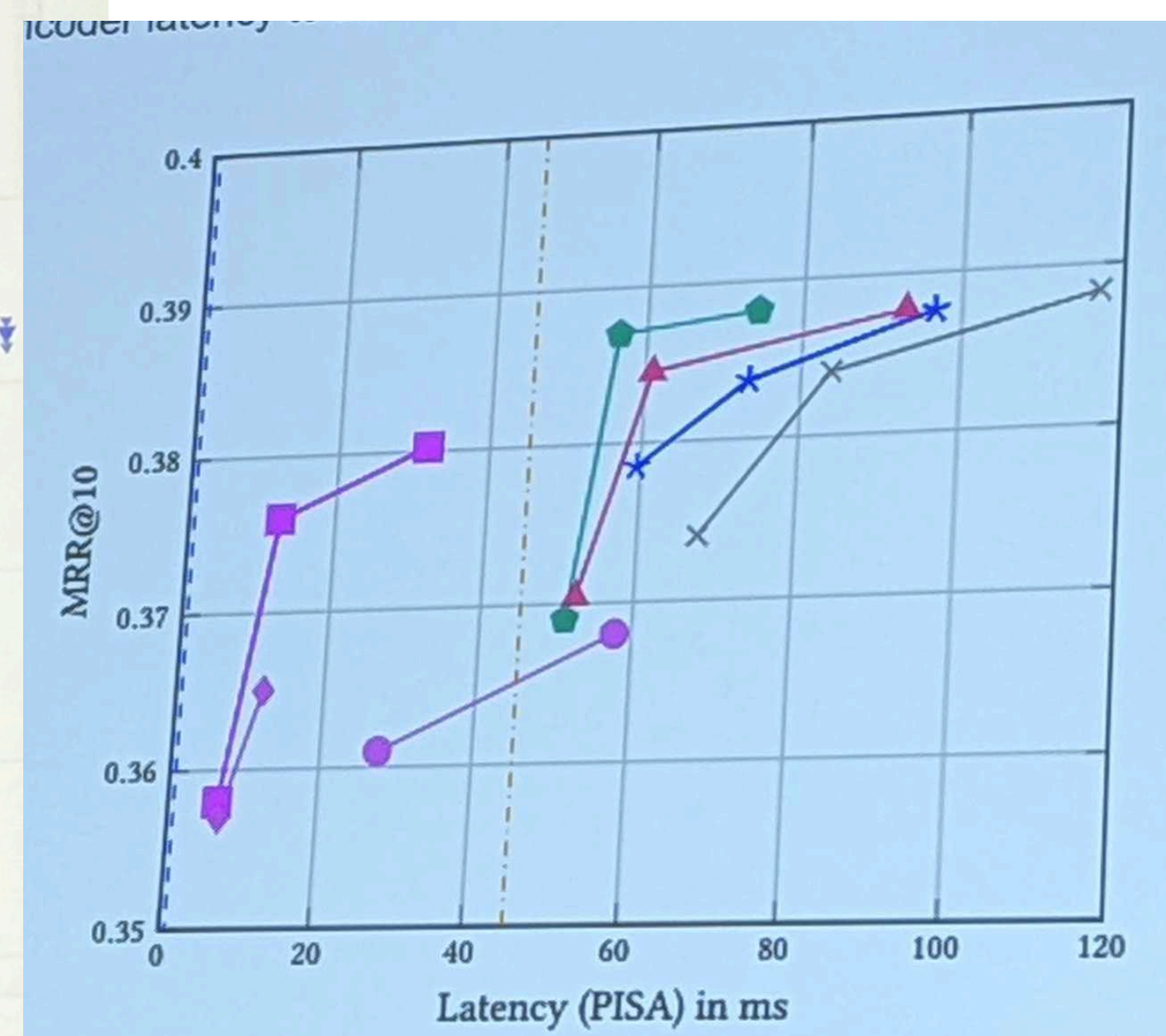
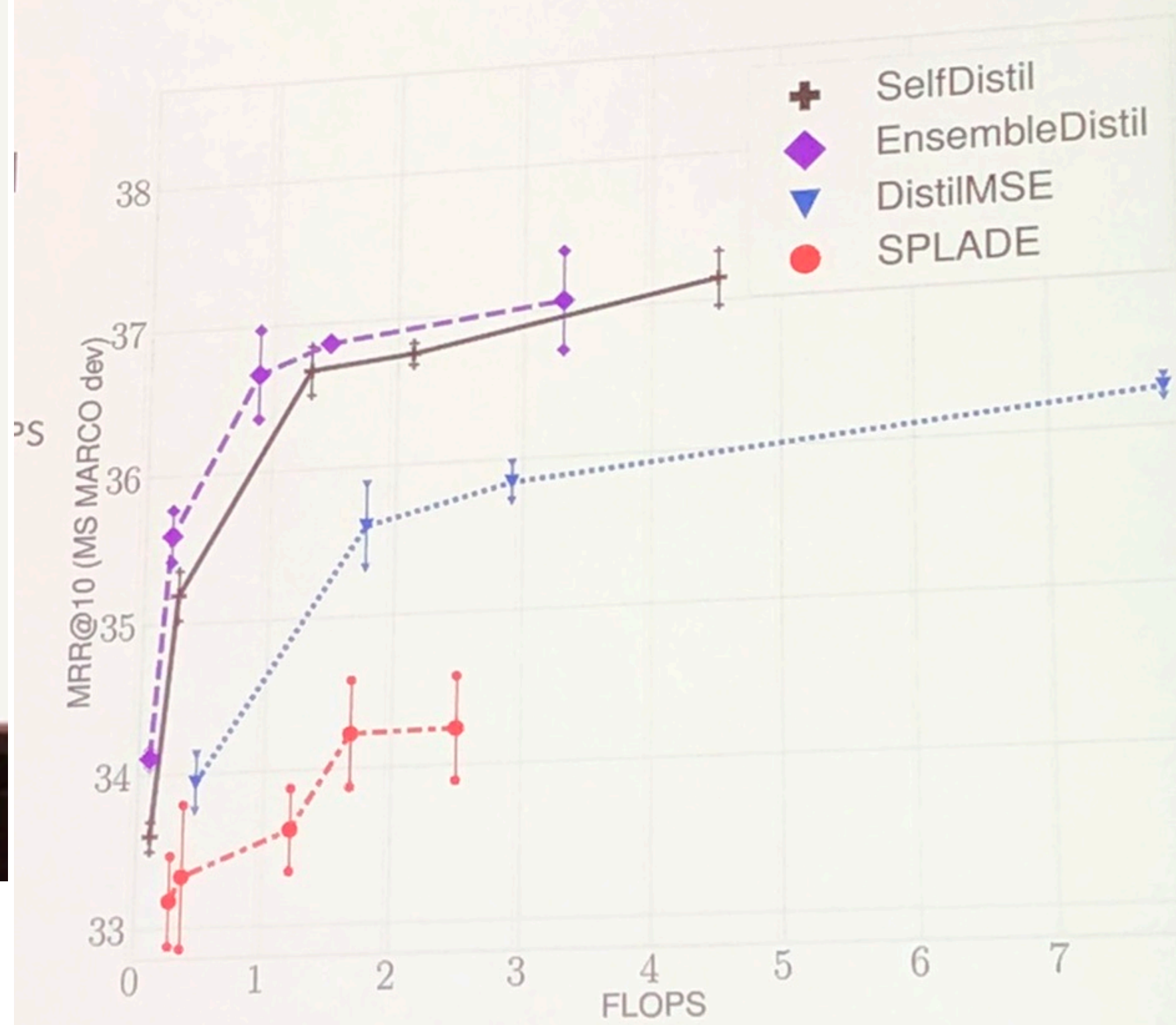
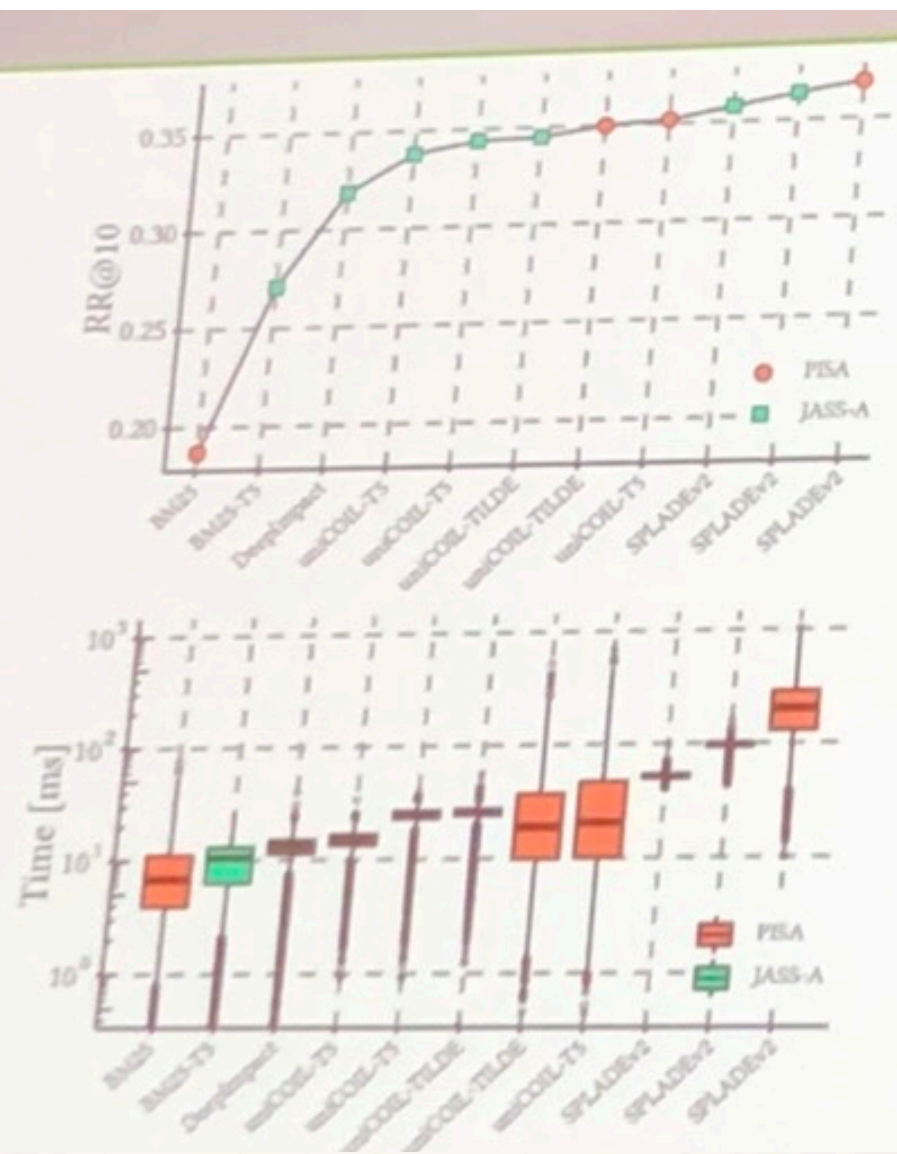


TREC Deep Learning Track 2019 – Document Ranking





TREC Deep Learning Track 2019 – Document Ranking



Developing Energy Efficient Filtering Systems

Leif Azzopardi, Wim Vanderbauwhede, Mahmoud Moadeli
Dept. of Comp. Sci., University of Glasgow
Glasgow, United Kingdom
{leif, wim, mahmoudm}@dcs.gla.ac.uk

ABSTRACT

Processing large volumes of information generally requires massive amounts of computational power, which consumes a significant amount of energy. An emerging challenge is the development of “environmentally friendly” systems that are not only efficient in terms of time, but also energy efficient. In this poster, we outline our initial efforts at developing greener filtering systems by employing Field Programmable Gate Arrays (FPGA) to perform the core information processing task. FPGAs enable code to be executed in parallel at a chip level, while consuming only a fraction of the power of a standard (von Neuman style) processor. On a number of test collections, we demonstrate that the FPGA filtering system performs 10-20 times faster than the Itanium based implementation, resulting in considerable energy savings.

Categories and Subject Descriptors

H.3.4 [Information Storage and Retrieval]: Systems and Software: Performance evaluation

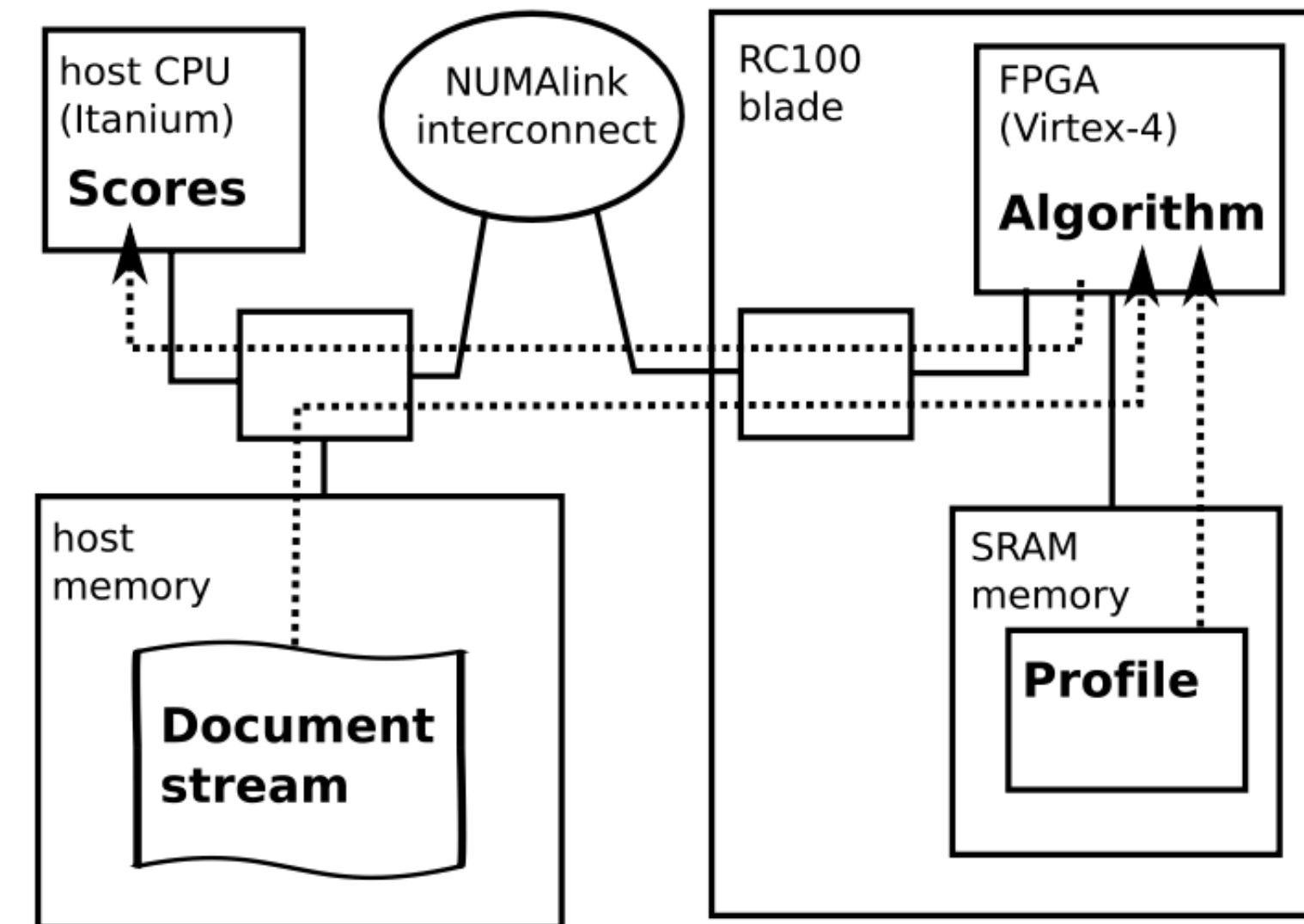
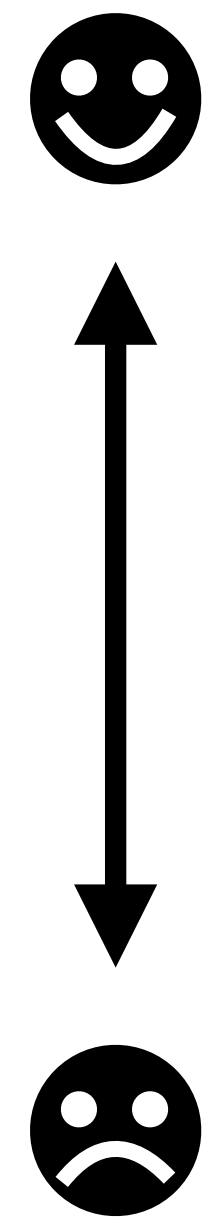


Figure 1: Schematic of FPGA-accelerated filtering

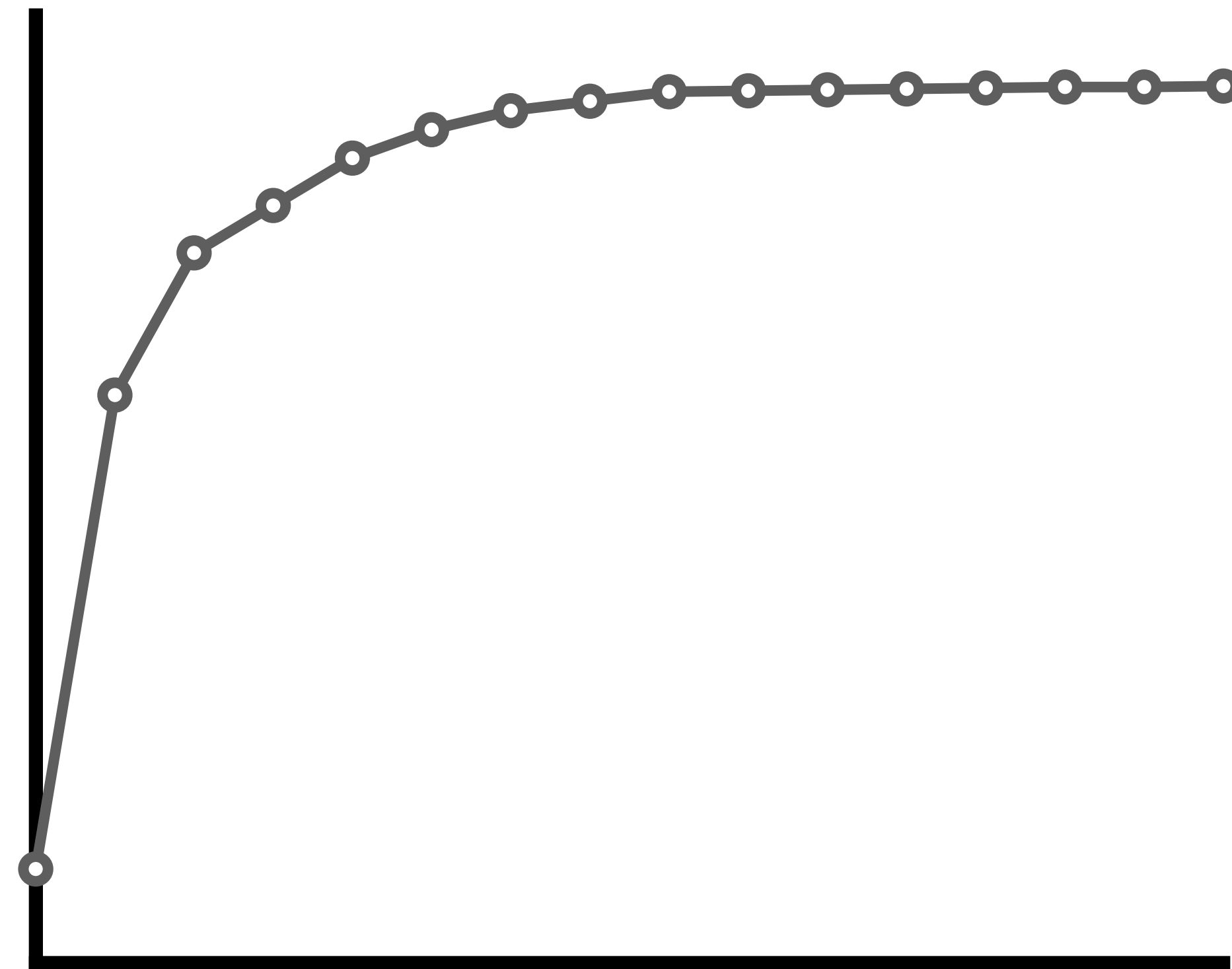
2. SYSTEM ARCHITECTURE

An FPGA is a reconfigurable semiconductor device which

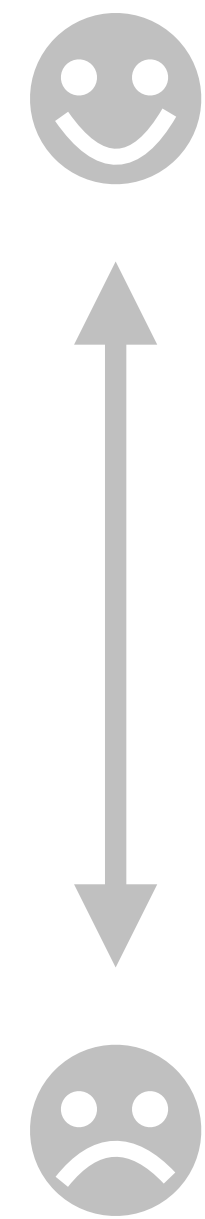




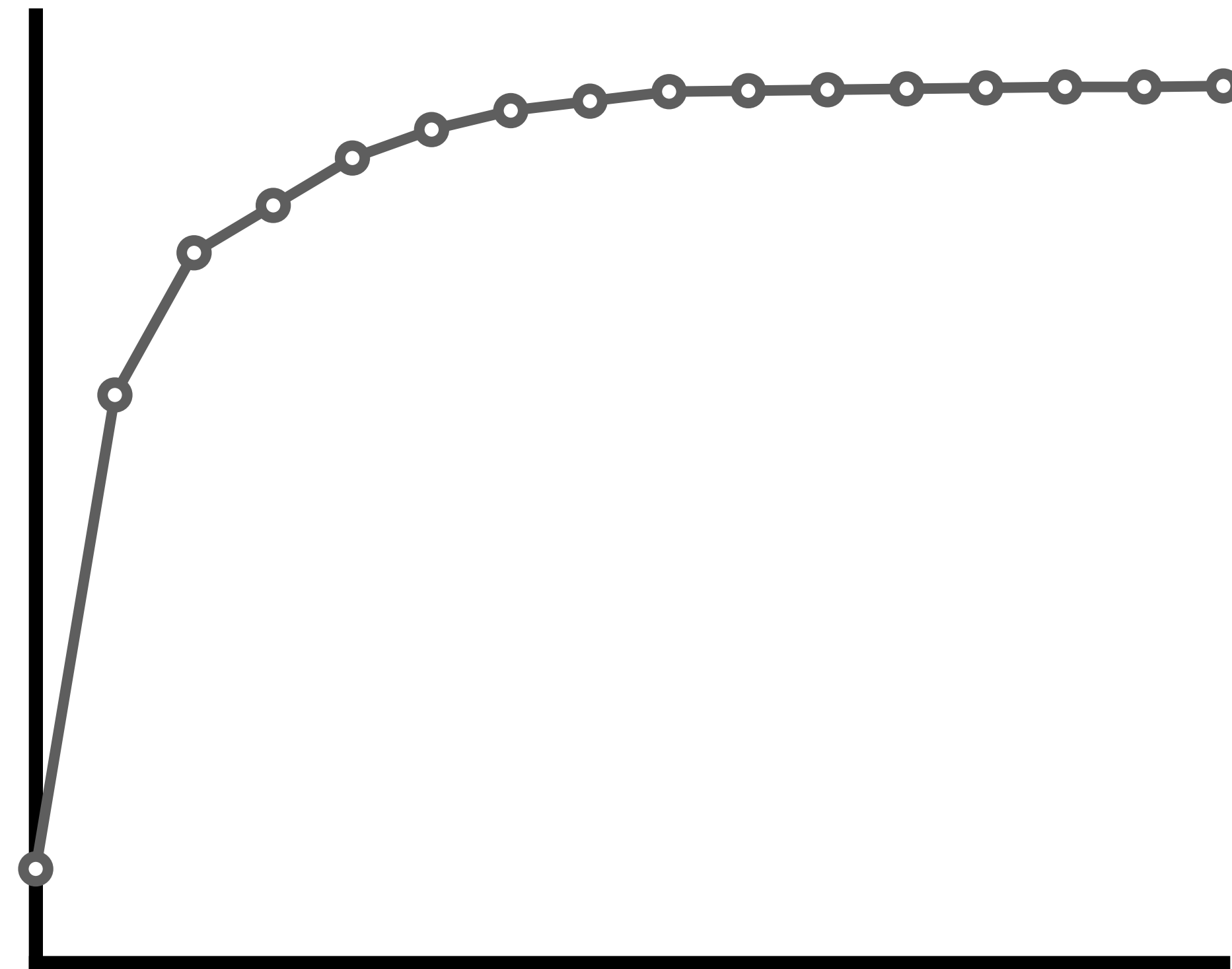
Effectiveness



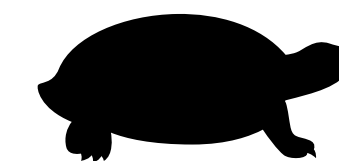
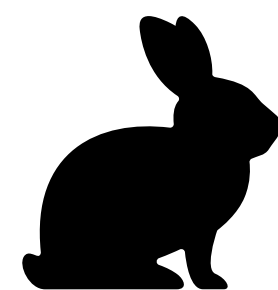
Efficiency

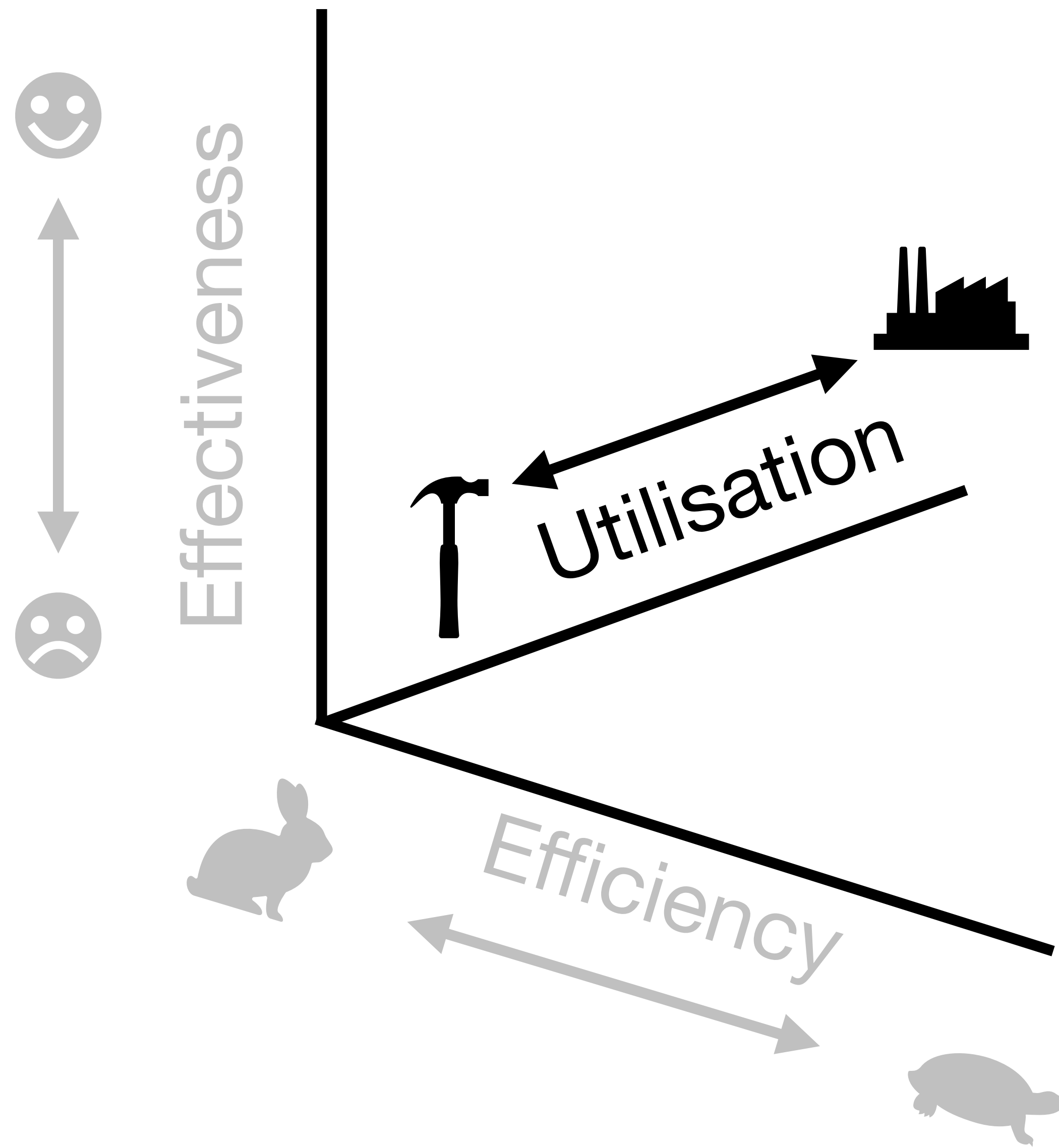


Effectiveness



Efficiency





Measuring emissions

- First, measure power consumption:

Measuring emissions

- First, measure power consumption:

$$p_t = \frac{\Omega \cdot t \cdot (p_c + p_r + p_g)}{1000}$$

Measuring emissions

- First, measure power consumption:

$$\text{watts} \rightarrow p_t = \frac{\Omega \cdot t \cdot (p_c + p_r + p_g)}{1000}$$

Measuring emissions

- First, measure power consumption:

$$\text{PUE} \rightarrow \Omega \cdot t \cdot (p_c + p_r + p_g)$$
$$\text{watts} \rightarrow p_t = \frac{\Omega \cdot t \cdot (p_c + p_r + p_g)}{1000}$$

Measuring emissions

- First, measure power consumption:

$$\text{watts} \rightarrow p_t = \frac{\text{PUE} \cdot \text{Running Time} \cdot (\underbrace{p_c + p_r + p_g}_{\text{watts}})}{1000}$$

Measuring emissions

- First, measure power consumption:

$$\text{watts} \rightarrow p_t = \frac{\text{PUE} \cdot \text{Running Time} \cdot \text{CPU, RAM, GPU power draw}}{1000}$$

The equation is annotated with arrows pointing to its components: 'PUE' points to the PUE term, 'Running Time' points to the t term, 'CPU, RAM, GPU power draw' points to the $(p_c + p_r + p_g)$ term, and 'watts' points to the p_t term on the left side of the equation.

Measuring emissions

- First, measure power consumption:

$$p_t = \frac{\text{PUE} \cdot \text{Running Time} \cdot (\text{CPU, RAM, GPU power draw})}{1000}$$

Diagram illustrating the formula for measuring power consumption:

The formula is: $p_t = \frac{\Omega \cdot t \cdot (p_c + p_r + p_g)}{1000}$

Annotations:

- Ω is labeled PUE (Power Usage Effectiveness).
- t is labeled Running Time.
- $(p_c + p_r + p_g)$ is labeled CPU, RAM, GPU power draw.
- p_t is labeled watts.

- Next, measure emissions:

Measuring emissions

- First, measure power consumption:

$$p_t = \frac{\text{PUE} \cdot \text{Running Time} \cdot (\text{CPU, RAM, GPU power draw})}{1000}$$

Diagram illustrating the formula for power consumption (p_t) in watts:

- PUE (Power Usage Effectiveness) is the first factor.
- Running Time (t) is the second factor.
- $(p_c + p_r + p_g)$ represents the power draw of the CPU, RAM, and GPU.
- The result is divided by 1000 to convert the value to watts.

- Next, measure emissions:

$$\text{kgCO}_2\text{e} = \theta \cdot p_t$$

Measuring emissions

- First, measure power consumption:

$$\text{watts} \rightarrow p_t = \frac{\text{PUE} \cdot \text{Running Time} \cdot (\text{CPU, RAM, GPU power draw})}{1000}$$

- Next, measure emissions:

$$\text{emissions} \rightarrow \text{kgCO}_2\text{e} = \theta \cdot p_t$$

Measuring emissions

- First, measure power consumption:

$$\text{watts} \rightarrow p_t = \frac{\text{PUE} \cdot \text{Running Time} \cdot (\text{CPU, RAM, GPU power draw})}{1000}$$

- Next, measure emissions:

$$\text{emissions} \rightarrow \text{kgCO}_2\text{e} = \theta \cdot p_t \leftarrow \text{Power consumption of experiments}$$

Measuring emissions

- First, measure power consumption:

$$\text{watts} \rightarrow p_t = \frac{\text{PUE} \cdot \text{Running Time} \cdot \text{CPU, RAM, GPU power draw}}{1000}$$

- Next, measure emissions:

$$\text{emissions} \rightarrow \text{kgCO}_2\text{e} = \theta \cdot p_t$$

avg. CO₂e (kg) per kWh where experiments took place

Power consumption of experiments

Measuring emissions

- First, measure power consumption:

$$\text{watts} \rightarrow p_t = \frac{\text{PUE} \cdot \text{Running Time} \cdot \text{CPU, RAM, GPU power draw}}{1000}$$

- Next, measure emissions:

$$\text{emissions} \rightarrow \text{kgCO}_2\text{e} = \theta \cdot p_t$$

avg. CO₂e (kg) per kWh where experiments took place

Power consumption of experiments

- Emissions of my search engine:

$$\Delta \text{kgCO}_2\text{e} = \theta \cdot \Delta_q \cdot p_q$$

Measuring emissions

- First, measure power consumption:

$$\text{watts} \rightarrow p_t = \frac{\text{PUE} \cdot \text{Running Time} \cdot \text{CPU, RAM, GPU power draw}}{1000}$$

- Next, measure emissions:

$$\text{emissions} \rightarrow \text{kgCO}_2\text{e} = \theta \cdot p_t$$

avg. CO₂e (kg) per kWh where experiments took place

Power consumption of experiments

- Emissions of my search engine:

$$\Delta \text{kgCO}_2\text{e} = \theta \cdot \Delta_q \cdot p_q$$

Power consumption of a single query

Measuring emissions

- First, measure power consumption:

$$\text{watts} \rightarrow p_t = \frac{\text{PUE} \cdot \text{Running Time} \cdot \text{CPU, RAM, GPU power draw}}{1000}$$

- Next, measure emissions:

$$\text{emissions} \rightarrow \text{kgCO}_2\text{e} = \theta \cdot p_t$$

avg. CO₂e (kg) per kWh where experiments took place

Power consumption of experiments

- Emissions of my search engine:

$$\Delta \text{kgCO}_2\text{e} = \theta \cdot \Delta_q \cdot p_q$$

No. queries issued per unit time

Power consumption of a single query

Measuring energy & emissions of your model

Name	CPU	DRAM	GPU	Network	Repository
CodeCarbon [71]	✓	✓	✓	✗	https://github.com/mlco2/codecarbon
pyJoules	✓	✓	✓	✗	https://github.com/powerapi-ng/pyJoules
energyusage [47]	✓	✓	✓	✗	https://github.com/responsibleproblemsolving/energy-usage
Carbontracker [3]	✓	✗	✓	✗	https://github.com/lflwa/carbontracker
Experiment Impact Tracker [33]	✓	✗	✓	✗	https://github.com/Breakend/experiment-impact-tracker
Cumulator [81]	✓	✓	✓	✓	https://github.com/epfl-iglobalhealth/cumulator

```
from codecarbon import EmissionsTracker
```

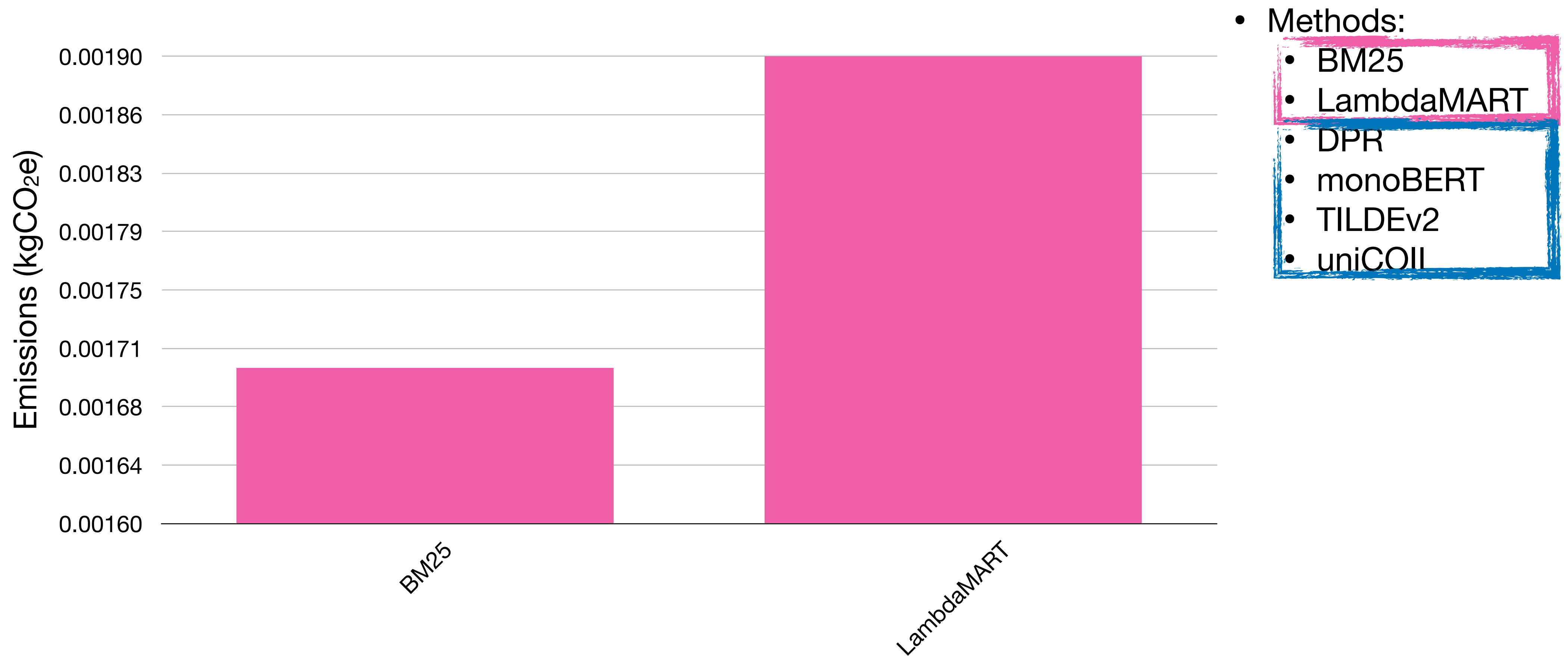
```
tracker = EmissionsTracker()
```

```
tracker.start()
```

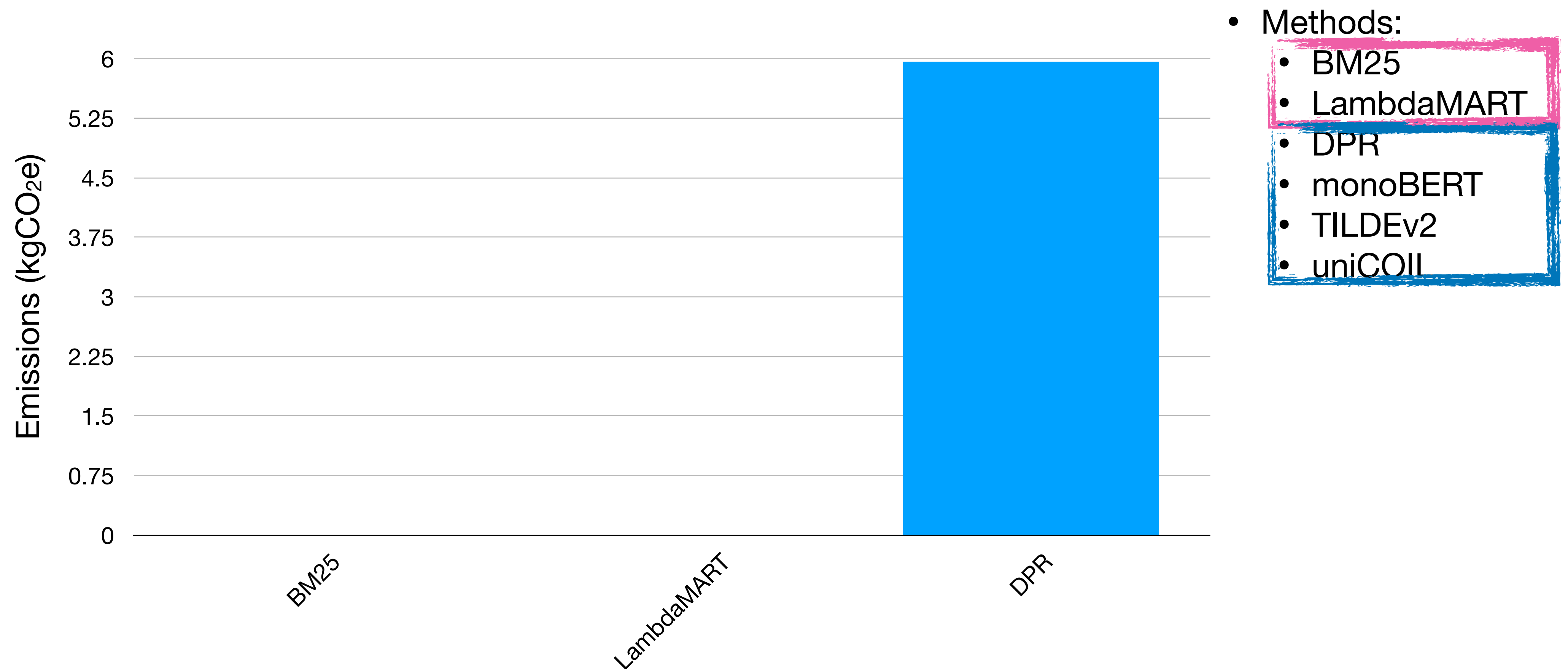
```
# Experiment code goes here
```

```
tracker.stop()
```

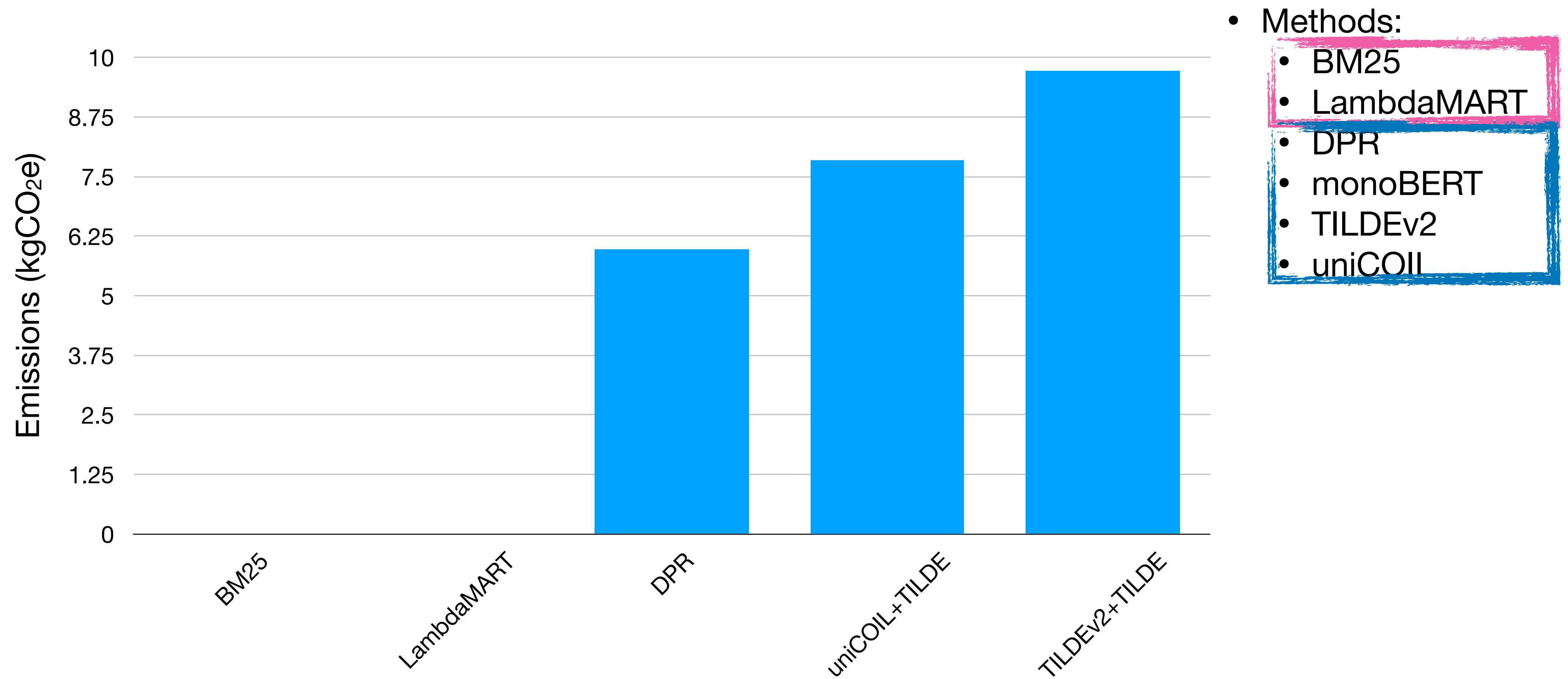
How many emissions do these methods produce to obtain an experimental result?



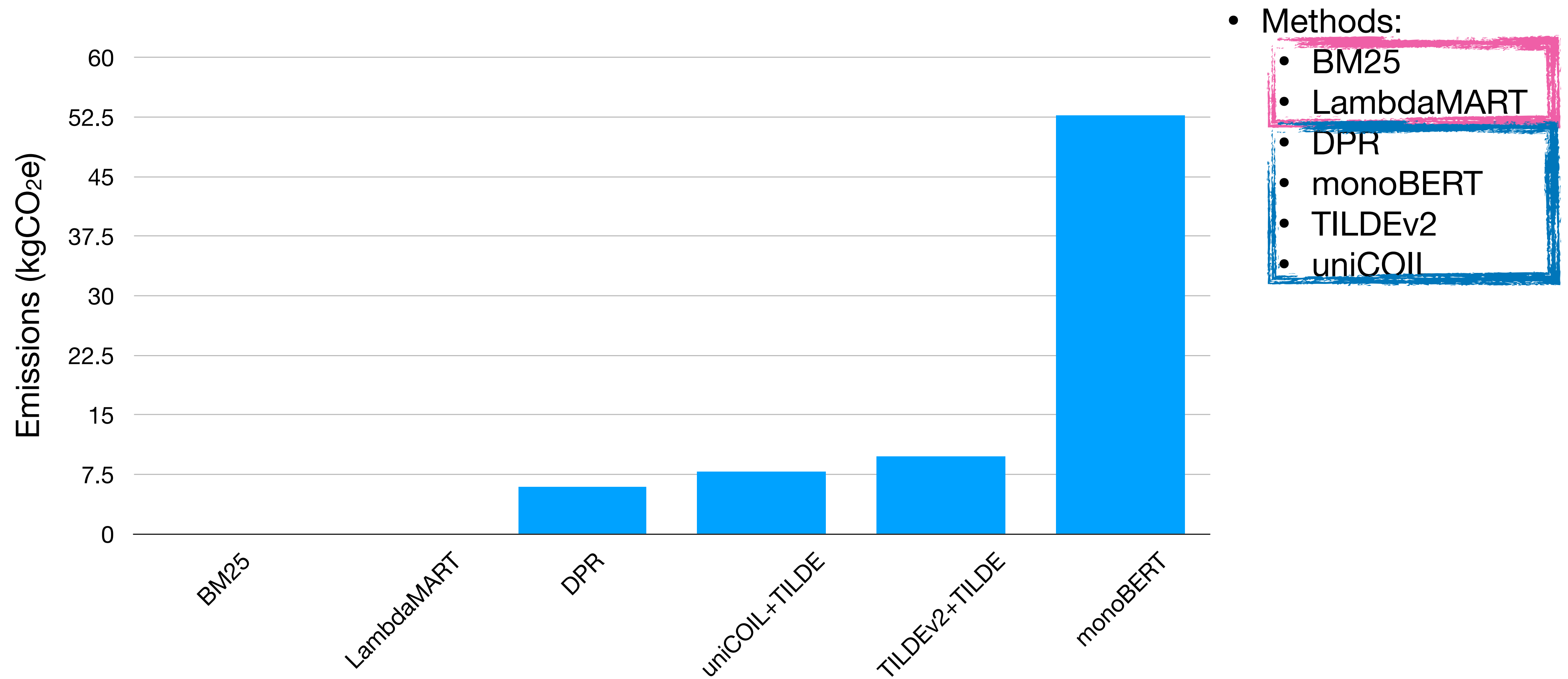
How many emissions do these methods produce to obtain an experimental result?



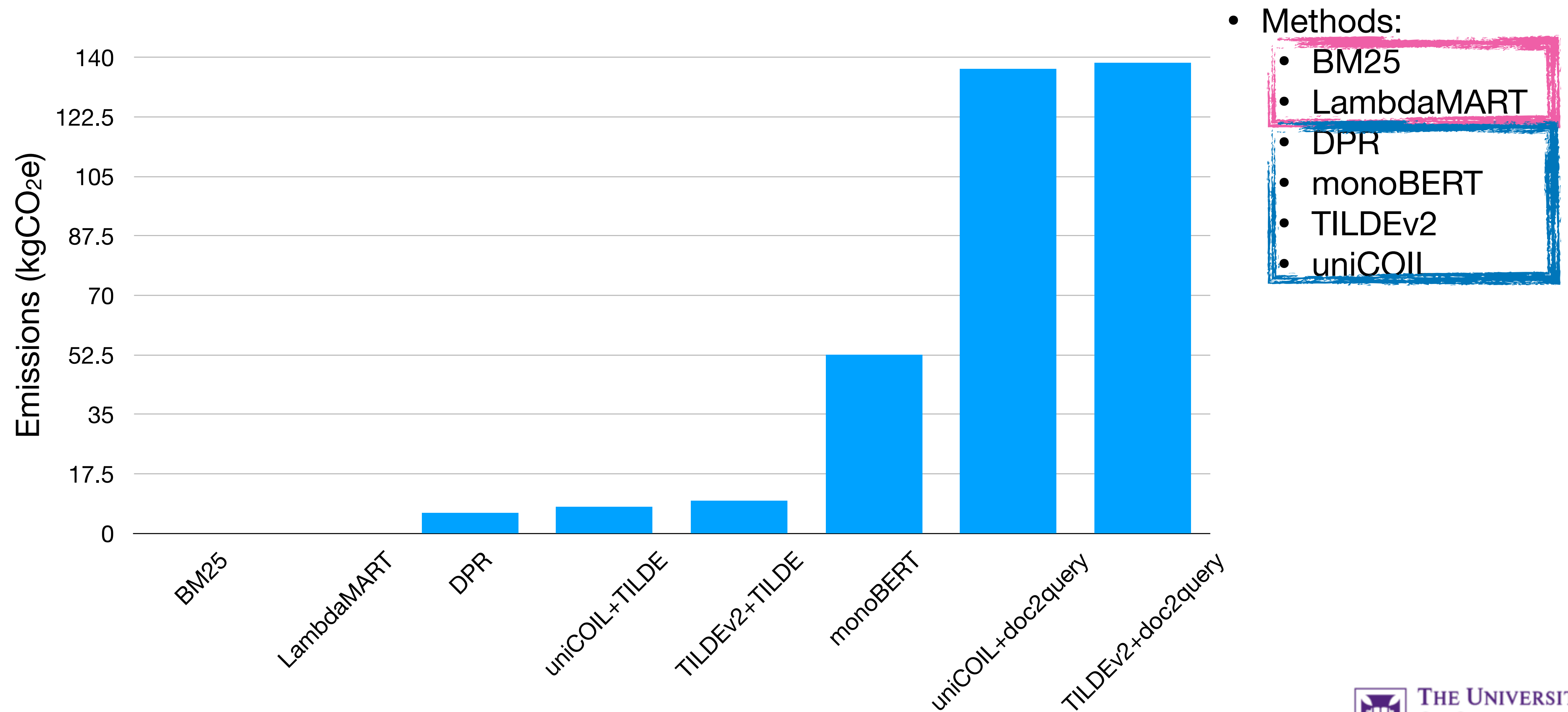
How many emissions do these methods produce to obtain an experimental result?



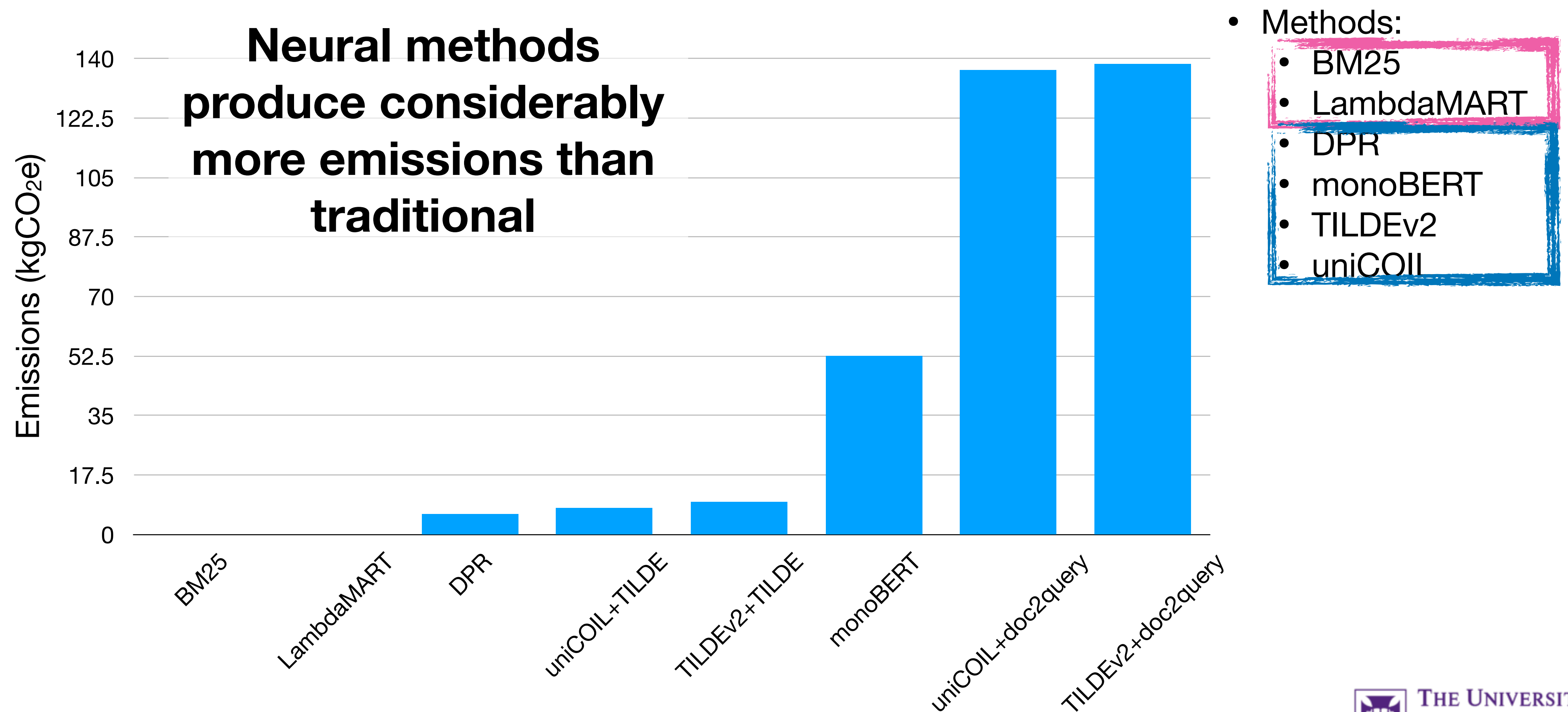
How many emissions do these methods produce to obtain an experimental result?



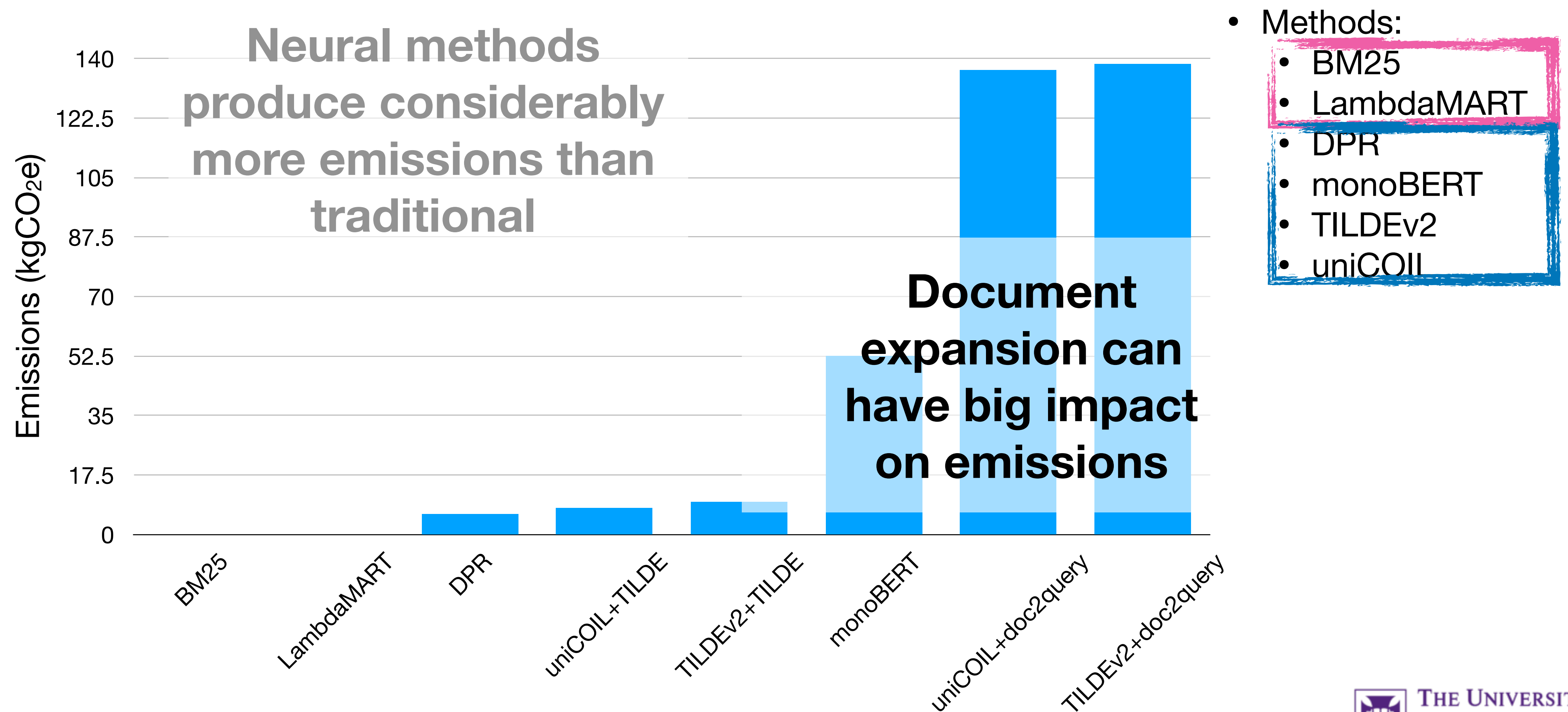
How many emissions do these methods produce to obtain an experimental result?



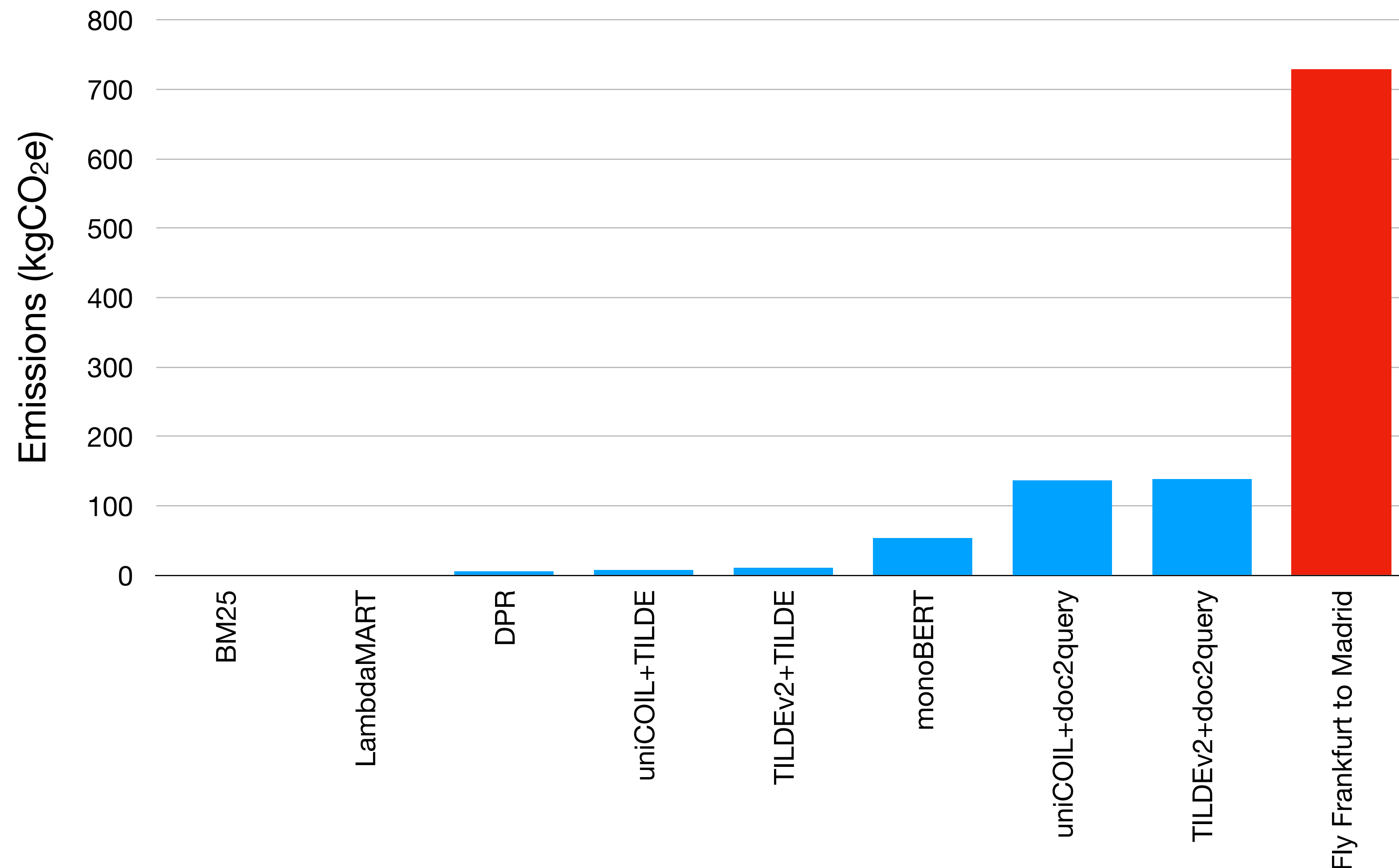
How many emissions do these methods produce to obtain an experimental result?



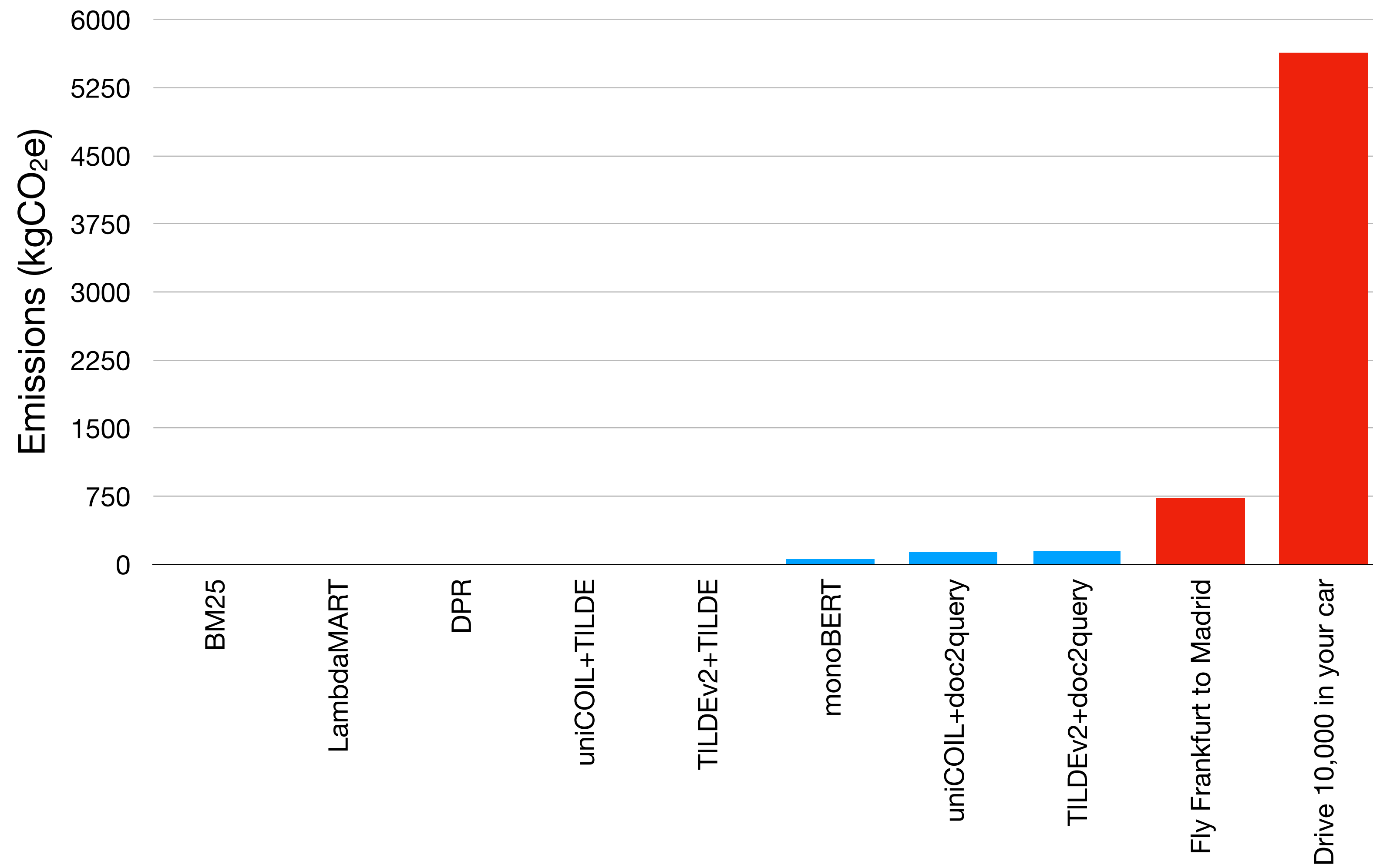
How many emissions do these methods produce to obtain an experimental result?



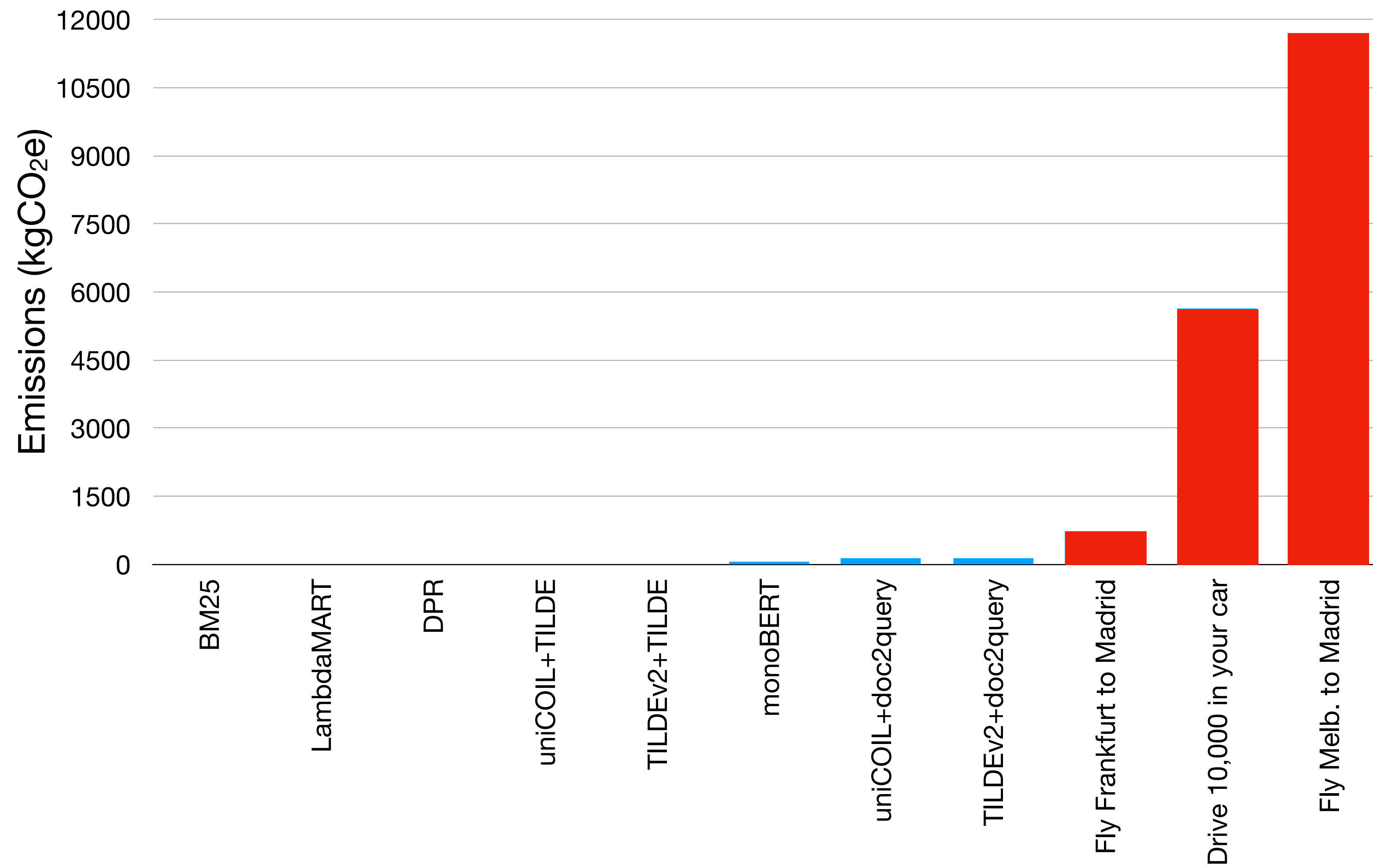
How many emissions do these methods produce to obtain an experimental result?



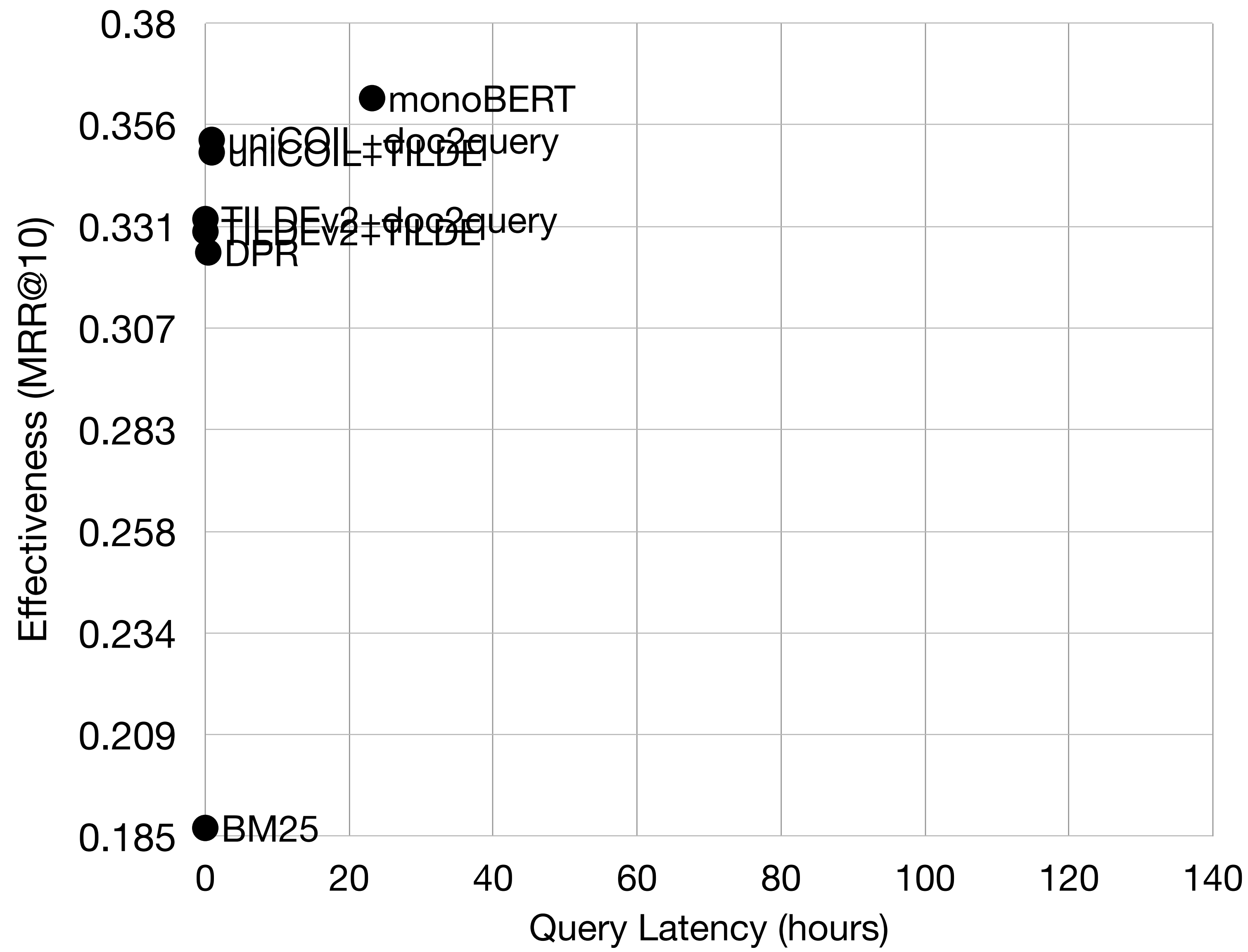
How many emissions do these methods produce to obtain an experimental result?



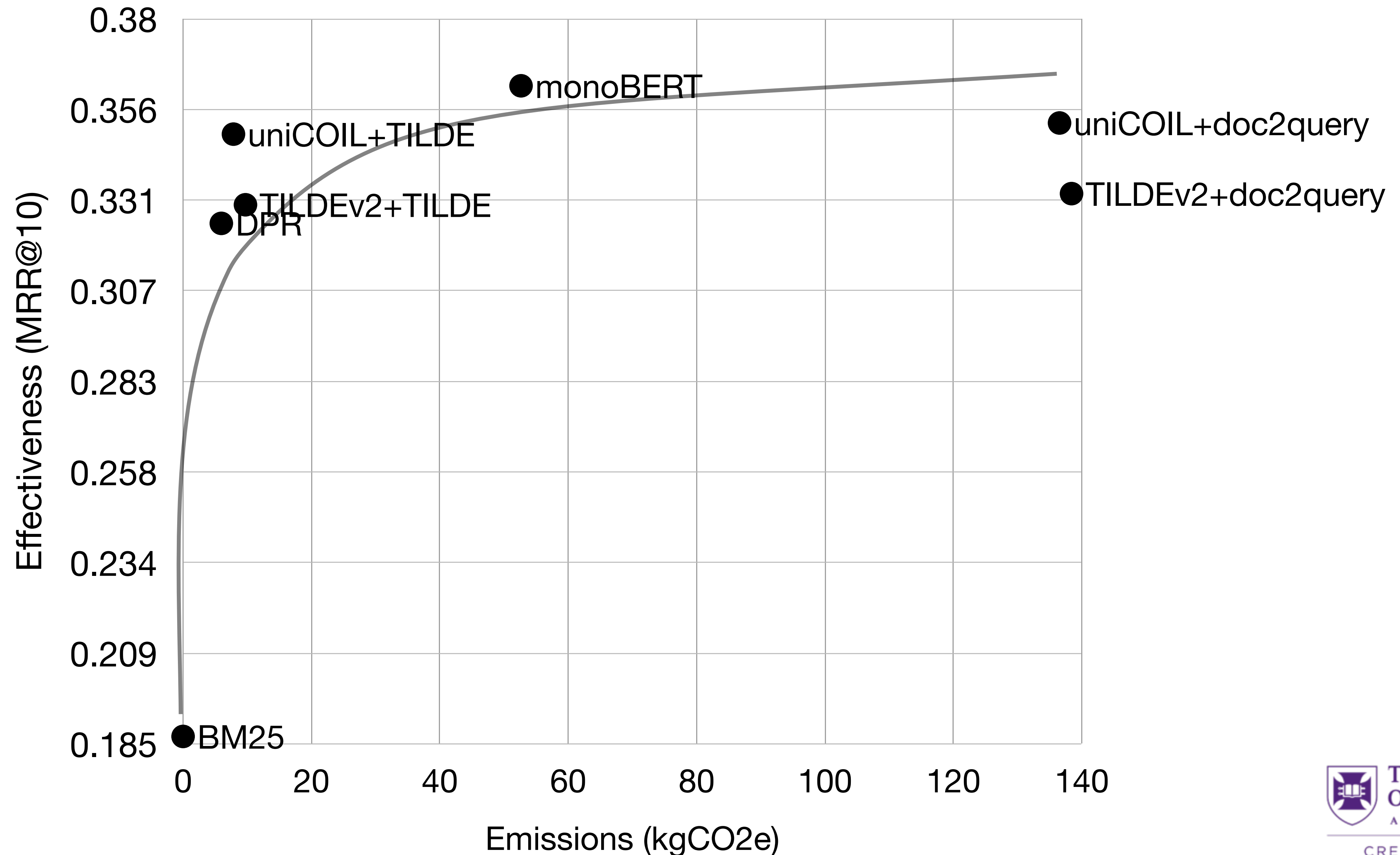
How many emissions do these methods produce to obtain an experimental result?



What are the effectiveness-utilisation trade-offs of these methods?



What are the effectiveness-utilisation trade-offs of these methods?







a framework for IR practitioners to remain mindful
of the potential costs of IR research

Reduce



VS



Reduce



VS

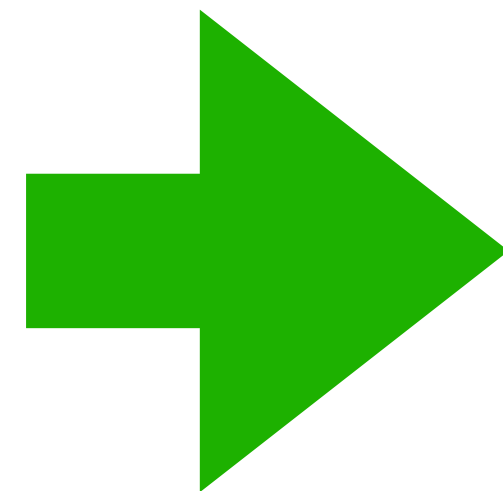


expend fewer resources

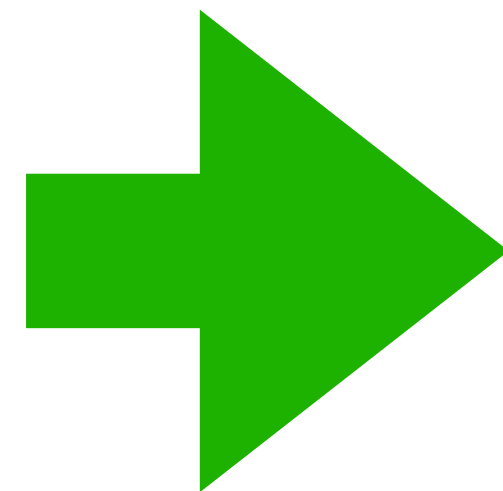
Reduce

- straightforward: simply reduce the number of experiments
- limit expensive computations, e.g., use CPU, FPGAs over GPU
- prior to starting any research or experiments, ask: *How can I perform research with fewer resources?*
 - Random Hyper-parameter Search
 - CPU-based Inference

Reuse



Reuse

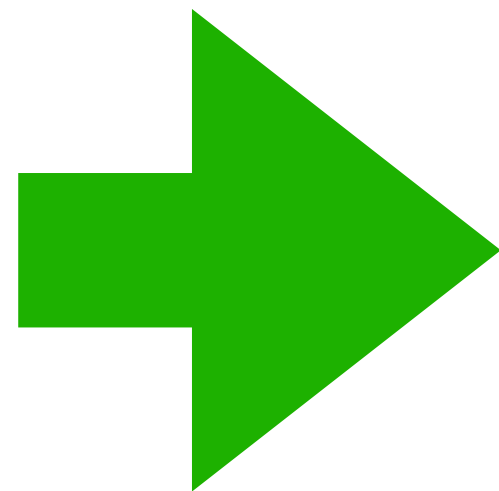


repurpose resources intended for one task to the same task

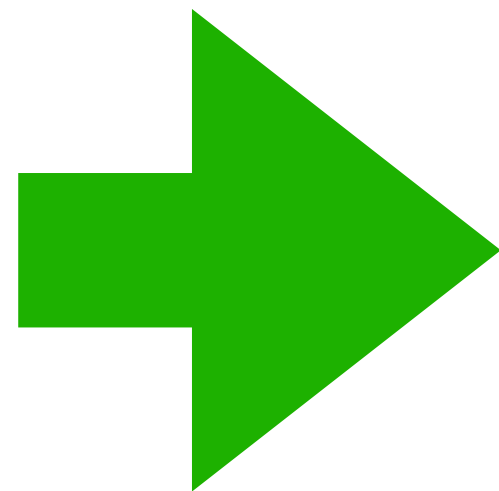
Reuse

- reuse existing software artefacts such as data, code, or models
- take something existing and repurpose it for the same task it was devised for
- prior to starting any research or experiments, ask: *How can I repurpose data, code, or other digital artefacts meant for one task to the same task?*
 - Reuse Large Collections
 - Pre-indexing Common Collections

Recycle



Recycle



repurpose resources intended for one task to a different task

Recycle

- recycle existing software artefacts such as data, code, or models
- recycle: the action of repurposing an existing artefact for a task it was not originally intended for
- prior to starting any research or experiments, ask: *How can I repurpose existing data, code, or other digital artefacts meant for one task to a different task?*
 - Neural Query Expansion
 - Passage expansion with TILDE

Outlook

- **Larger neural methods** = power-hungry hardware = utilisation of more power
 - but: increase model size for higher effectiveness may not apply to IR, as it does to NLP and ML

Outlook

- **Larger neural methods** = power-hungry hardware = utilisation of more power
 - but: increase model size for higher effectiveness may not apply to IR, as it does to NLP and ML
- Likely trend in neural IR: go beyond PLMs designed for NLP but are specialised for IR -> **pre-train for IR**
 - more power and more emissions
 - DSI: end-to-end transformers that encapsulate entire indexing & searching architecture into single model

Outlook

- **Larger neural methods** = power-hungry hardware = utilisation of more power
 - but: increase model size for higher effectiveness may not apply to IR, as it does to NLP and ML
- Likely trend in neural IR: go beyond PLMs designed for NLP but are specialised for IR -> **pre-train for IR**
 - more power and more emissions
 - DSI: end-to-end transformers that encapsulate entire indexing & searching architecture into single model
- IR community at a **turning point**
 - Bigger/more complex models
 - Bigger collections

Outlook

- **Larger neural methods** = power-hungry hardware = utilisation of more power
 - but: increase model size for higher effectiveness may not apply to IR, as it does to NLP and ML
- Likely trend in neural IR: go beyond PLMs designed for NLP but are specialised for IR -> **pre-train for IR**
 - more power and more emissions
 - DSI: end-to-end transformers that encapsulate entire indexing & searching architecture into single model
- IR community at a **turning point**
 - Bigger/more complex models
 - Bigger collections
- Let's be mindful of the **cost** of IR research
 - Power usage → \$\$\$
 - Emissions → CO₂e