

# Green Information Retrieval Research

---

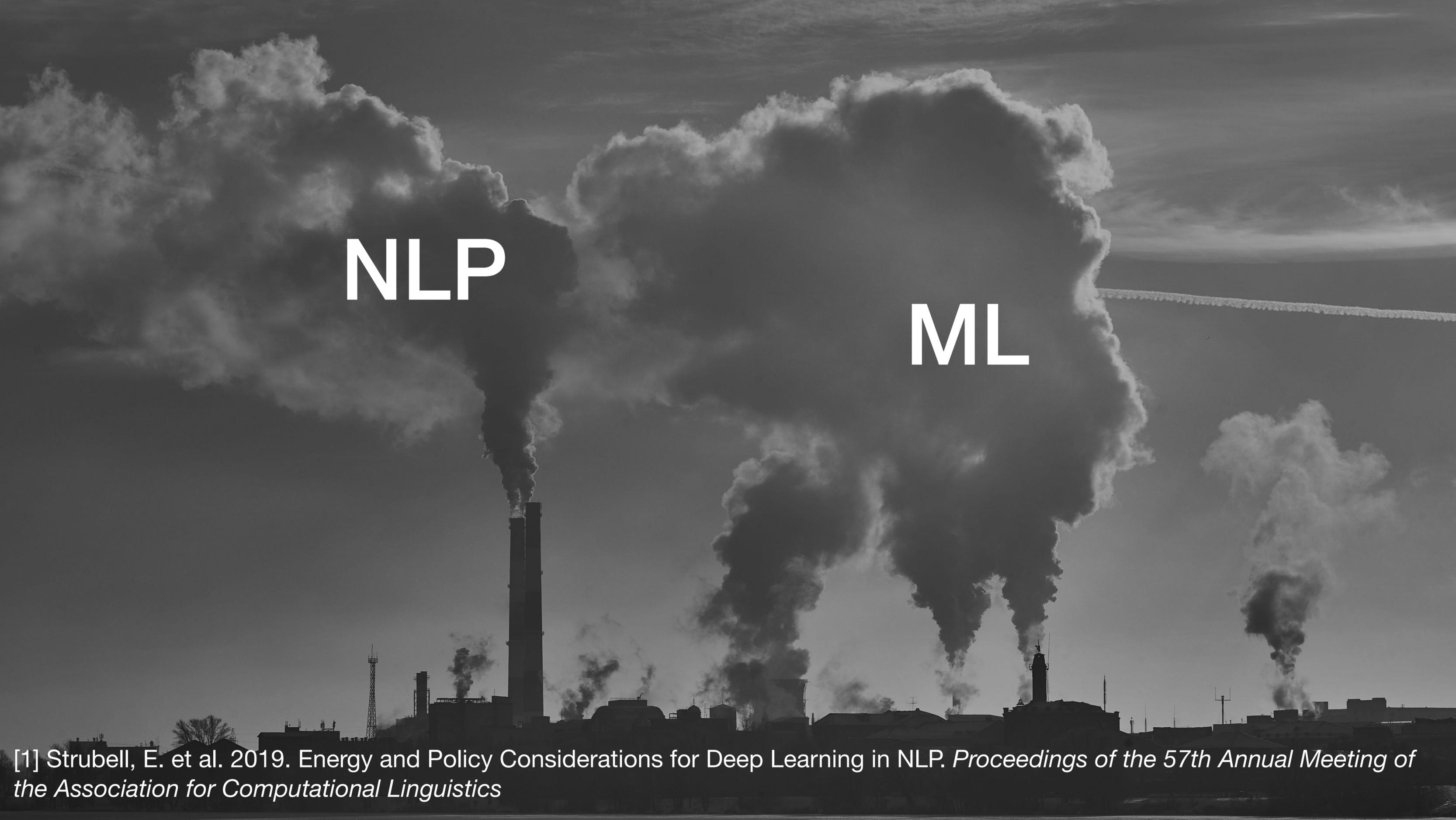


**Harry Scells**

Leipzig University, Germany

# PART I

*Context*



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*

# Why?

- Large (pre-trained) neural language models
  - Expend high energy for training and inference (compared to traditional models)
- The energy demands expected to continue growing as size and complexity of models increase
- Data centers and other infrastructure used to run these models also consume energy





NLP

ML

What about IR research?

# 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 kgCO<sub>2</sub>e; kilograms of carbon dioxide equivalent



NLP

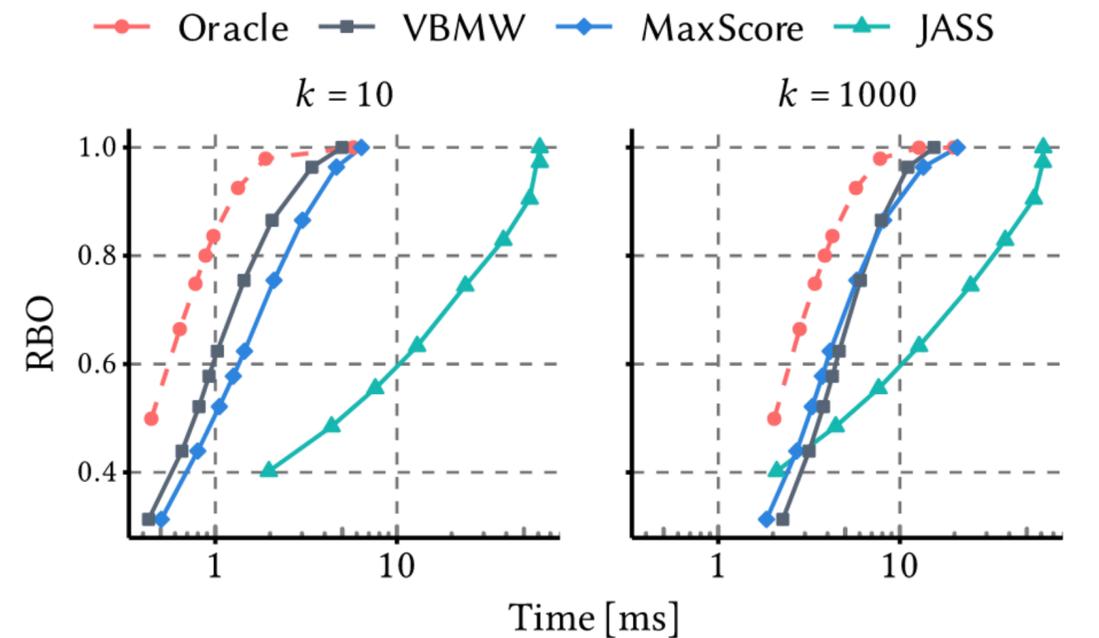
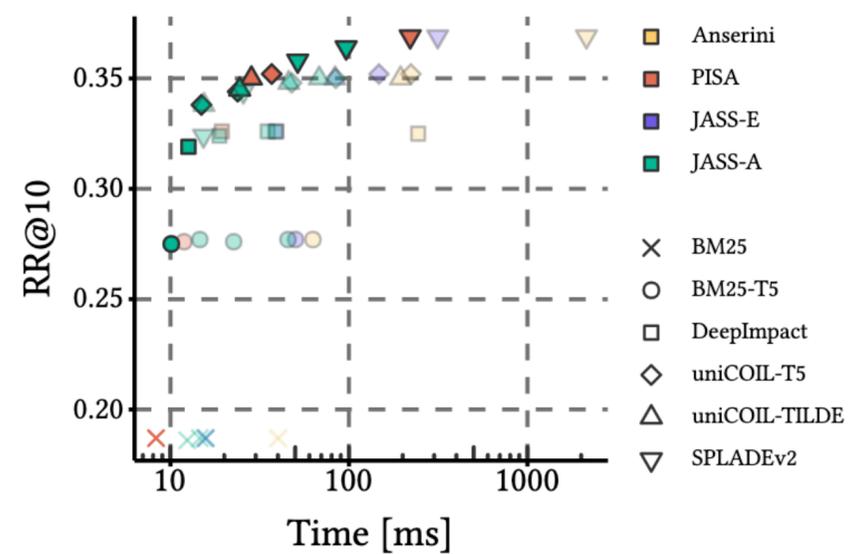
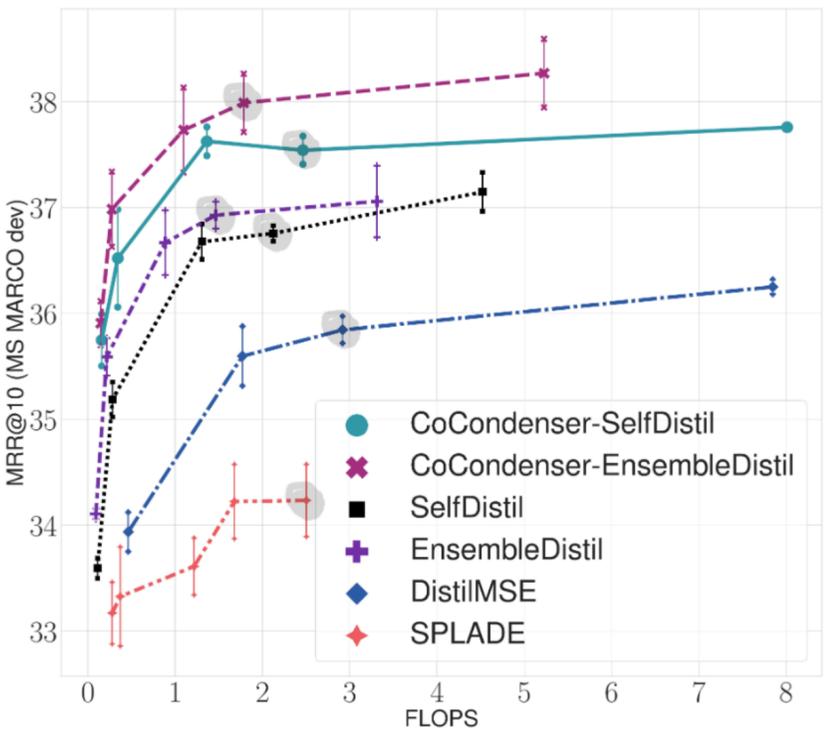
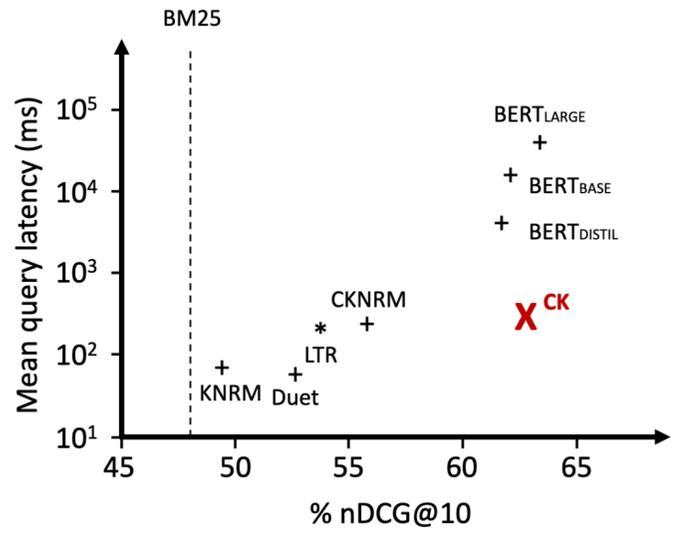
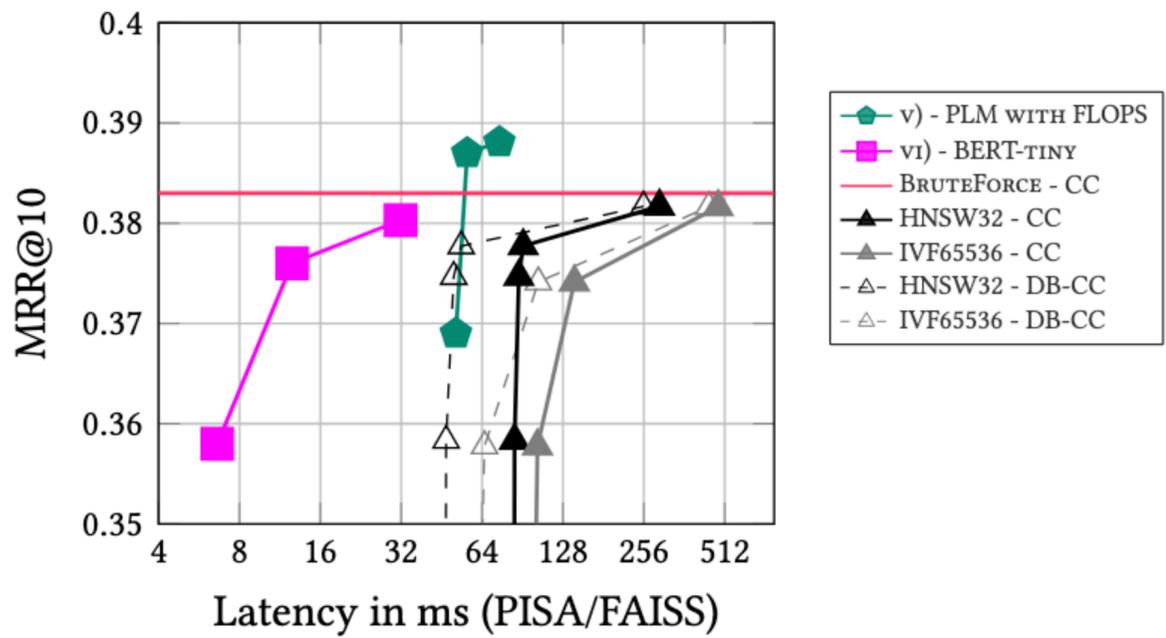
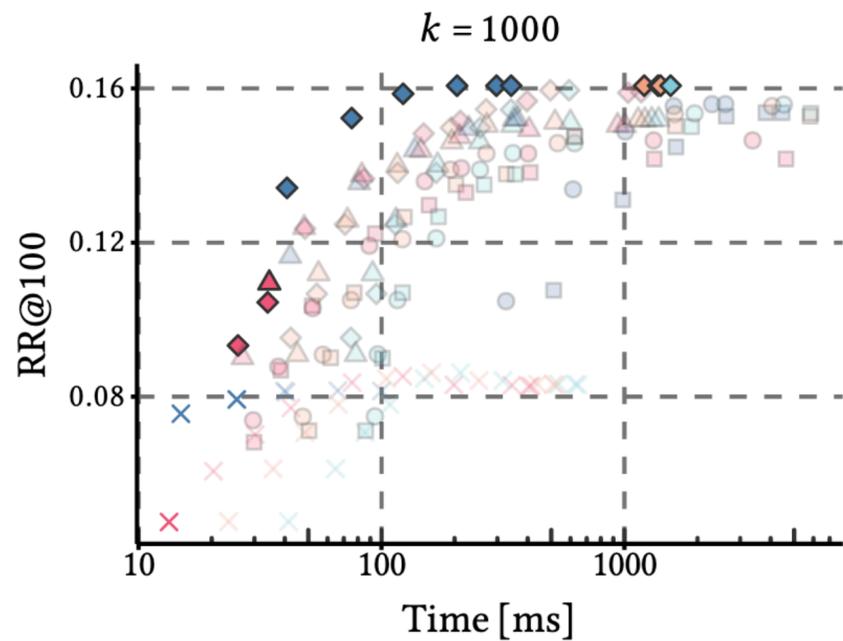
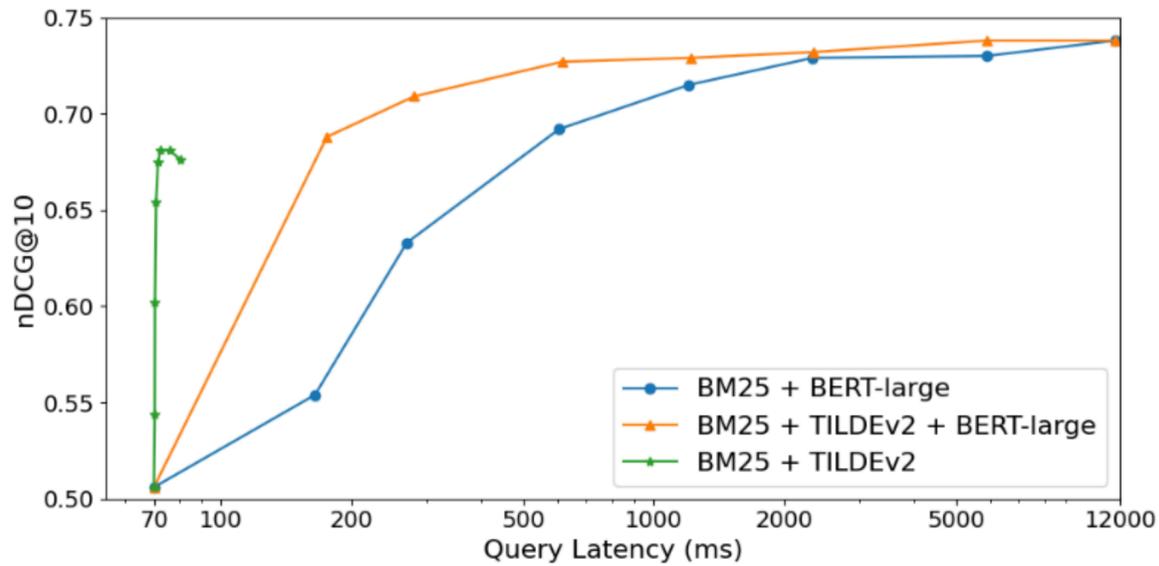
ML

What about IR research?

Isn't this just retrieval efficiency?

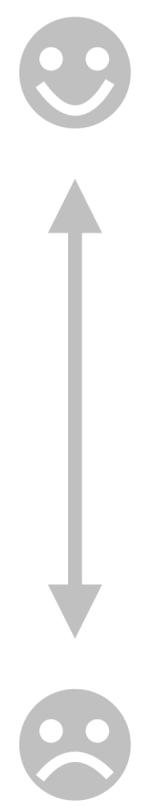
# Retrieval Efficiency

- **Speed** a system is able to retrieve relevant documents or information in response to a query.
- Factors that can impact retrieval efficiency include:
  - **Size and complexity of the corpus** being searched
  - Effectiveness of the **retrieval models** or techniques being used
  - Efficiency of the **hardware and infrastructure** used

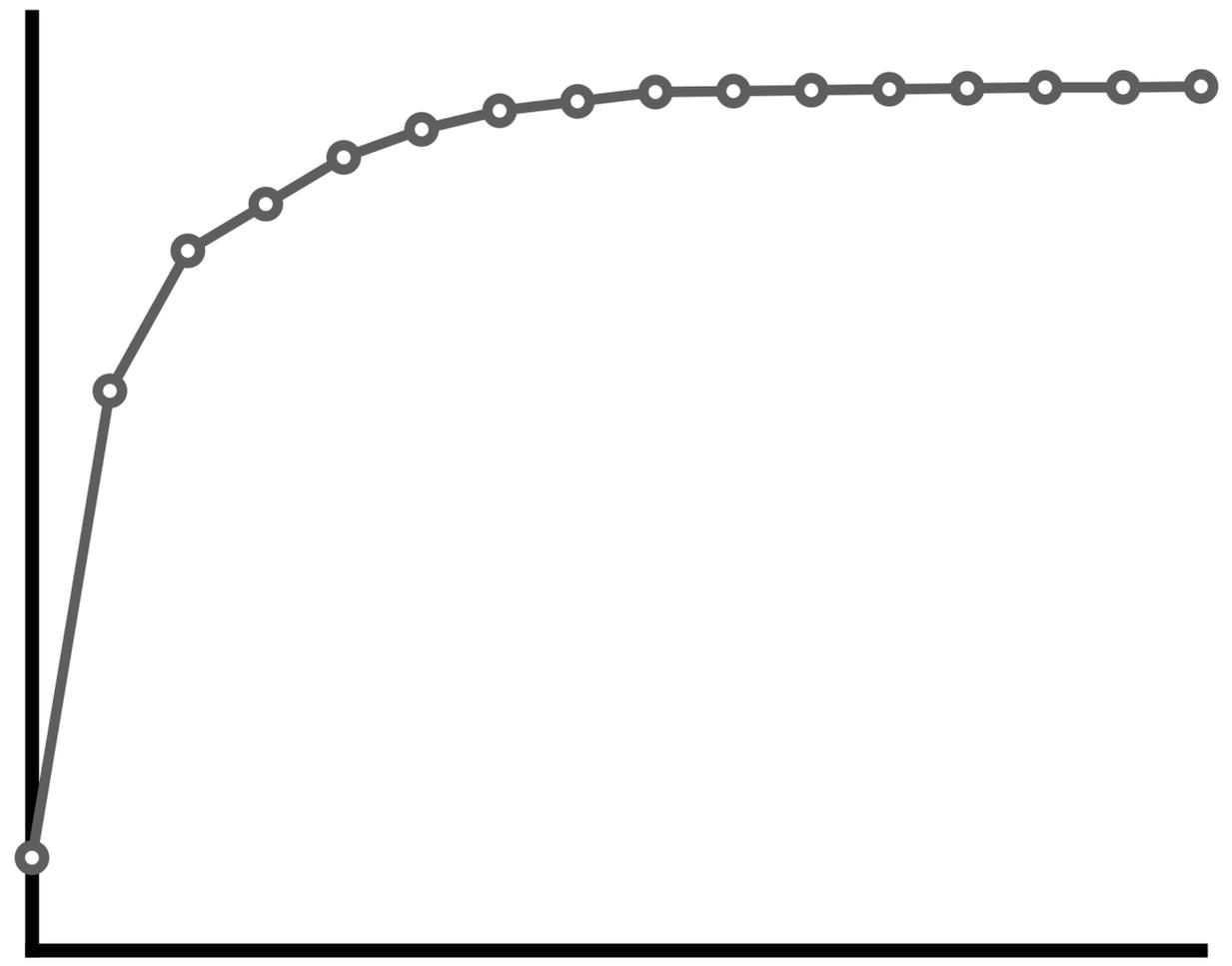






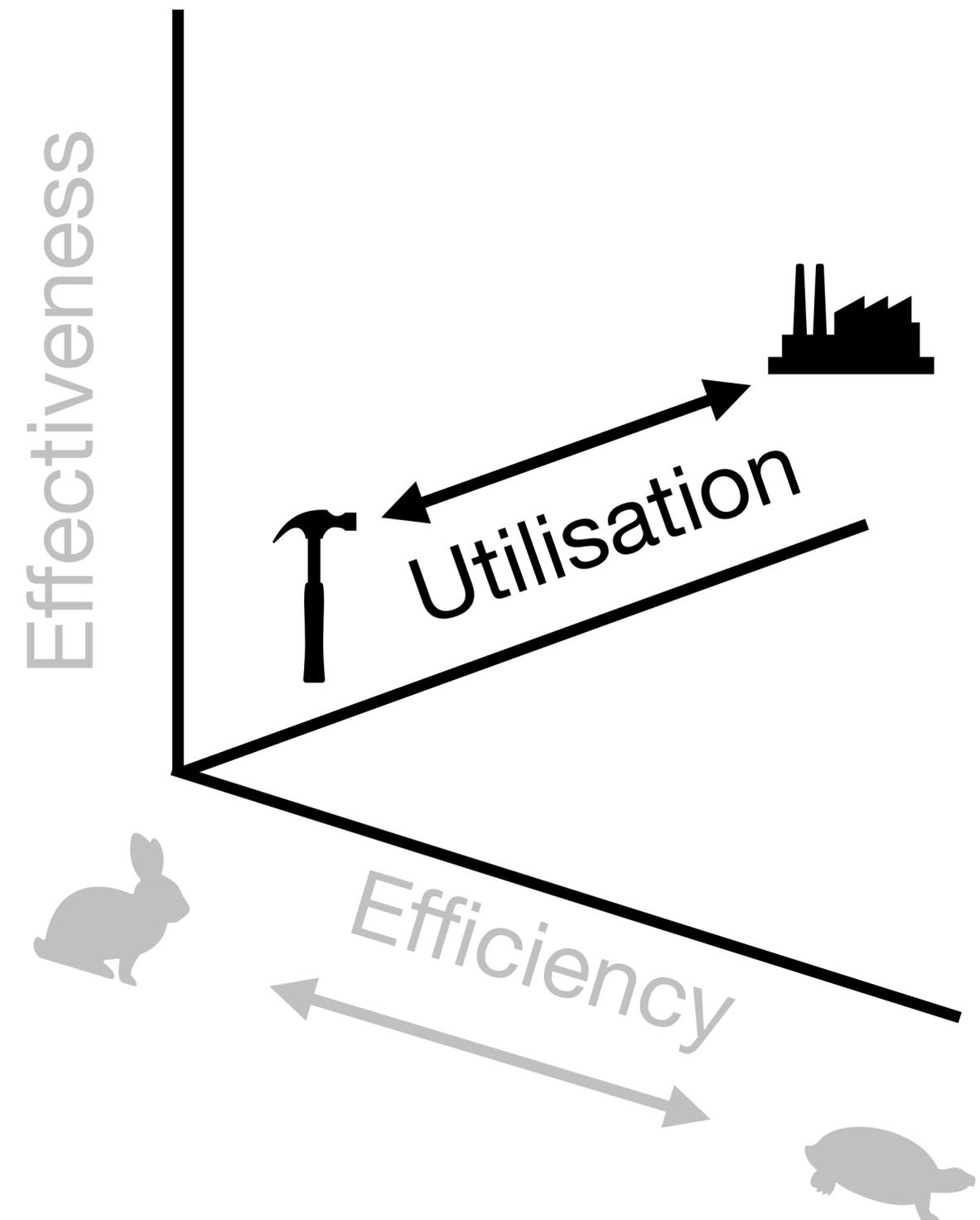


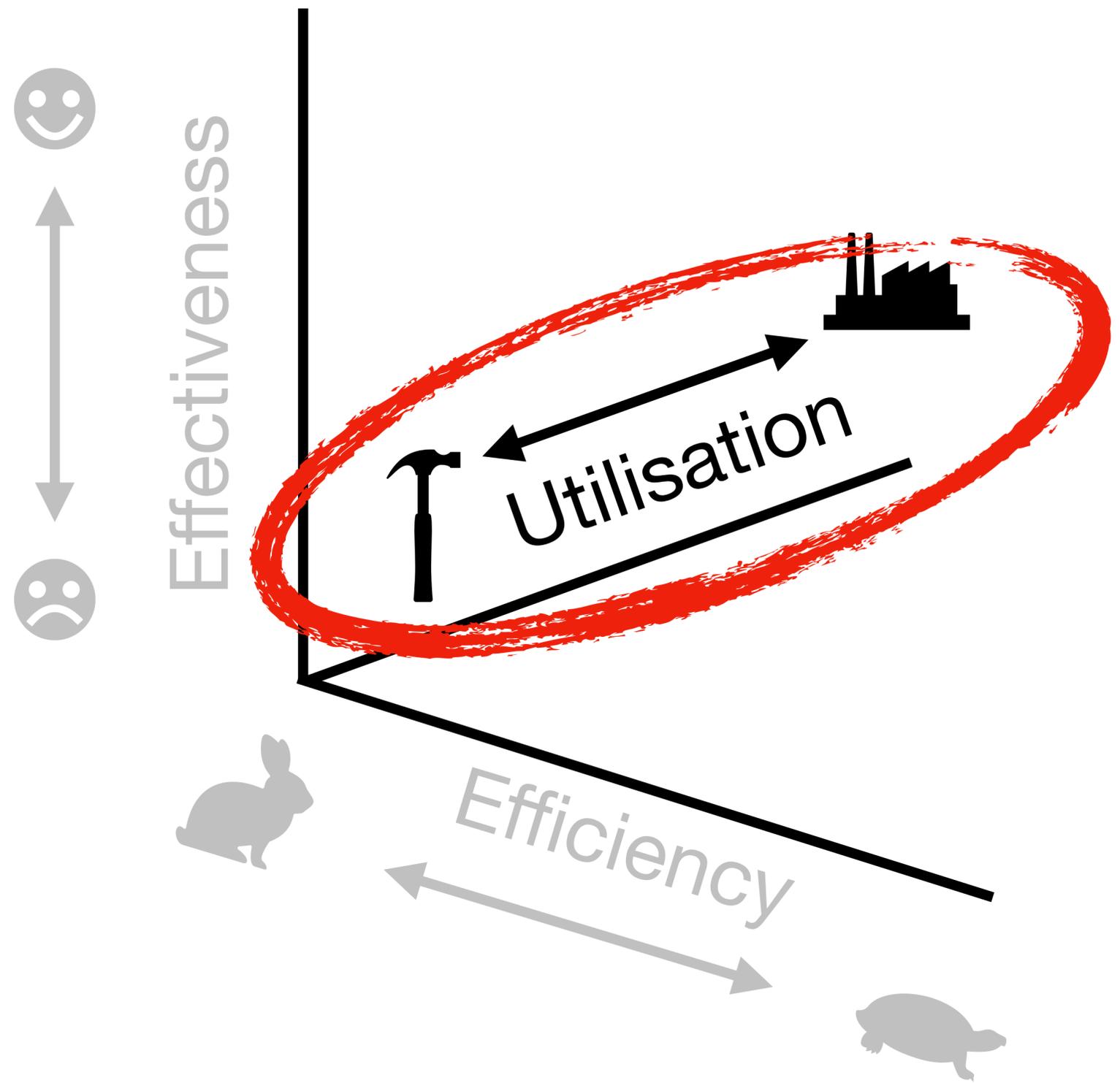
Effectiveness



Efficiency







**Okay, so what does  
this mean for IR?**

# Utilisation and Green IR

## **Green IR is...**

- *“research that yields novel results while taking into account the computational cost, encouraging a reduction in resources spent” [2]*

# Utilisation and Green IR

## Green IR is...

- “*research that yields novel results while taking into account the computational cost, encouraging a reduction in resources spent*” [2]
- Neural methods require pre-trained LMs
  - **Expensive** to create
  - Trend in IR towards creating **IR-specific** LMs [3,4,5,6]

[3] Gao, L. and Callan, J. 2021. Condenser: a Pre-training Architecture for Dense Retrieval. Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing

[4] Ma, X. et al. 2021. PROP: Pre-training with Representative Words Prediction for Ad-hoc Retrieval. Proceedings of the 14th ACM International Conference on Web Search and Data Mining

[5] Tay, Y. et al. 2022. Transformer Memory as a Differentiable Search Index. arXiv preprint arXiv:2202.06991.

[6] Zhou, Y. et al. 2022. DynamicRetriever: A Pre-training Model-based IR System with Neither Sparse nor Dense Index. arXiv preprint

[2] Schwartz, R. et al. 2020. Green AI. Communications of the ACM.

# Utilisation and Green IR

## Green IR is...

- “*research that yields novel results while taking into account the computational cost, encouraging a reduction in resources spent*” [2]
- Neural methods require pre-trained LMs
  - **Expensive** to create
  - Trend in IR towards creating **IR-specific LMs** [3,4,5,6]

**Pre-trained LMs come  
at a high power and  
emissions cost**

[3] Gao, L. and Callan, J. 2021. Condenser: a Pre-training Architecture for Dense Retrieval. Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing

[4] Ma, X. et al. 2021. PROP: Pre-training with Representative Words Prediction for Ad-hoc Retrieval. Proceedings of the 14th ACM International Conference on Web Search and Data Mining

[5] Tay, Y. et al. 2022. Transformer Memory as a Differentiable Search Index. arXiv preprint arXiv:2202.06991.

[6] Zhou, Y. et al. 2022. DynamicRetriever: A Pre-training Model-based IR System with Neither Sparse nor Dense Index. arXiv preprint

[2] Schwartz, R. et al. 2020. Green AI. Communications of the ACM.

# Utilisation and Green IR

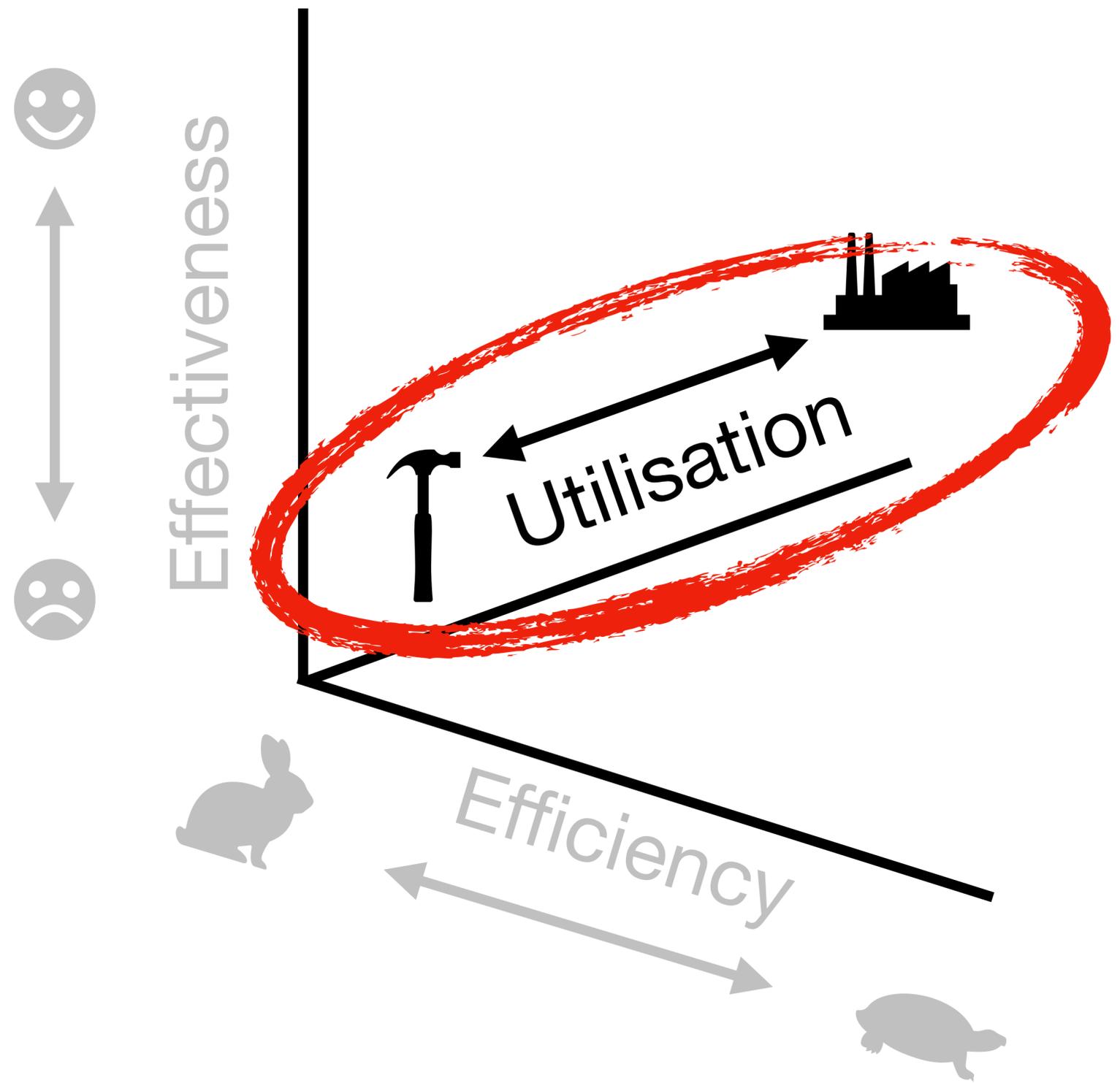
## Green IR is...

- “*research that yields novel results while taking into account the computational cost, encouraging a reduction in resources spent*” [2]

- Neural methods require pre-trained LMs
  - **Expensive** to create
  - Trend in IR towards creating **IR-specific LMs** [3,4,5,6]

**Pre-trained LMs come at a high power and emissions cost**

- Missing dimension of IR evaluation
  - Effectiveness
  - Efficiency
  - **Utilisation**



~~Okay, so what does this mean for IR?~~

**Okay, so how can I measure this?**

# 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:

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

# Measuring emissions

- First, measure power consumption:

$$\text{watts} \rightarrow p_t = \frac{\text{PUE} \cdot \text{Running Time} \cdot (\rho_c + \rho_r + \rho_g)}{1000}$$

# Measuring emissions

- First, measure power consumption:

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

The equation is annotated with arrows pointing to its components: 'watts' points to  $p_t$ , 'PUE' points to the PUE term, 'Running Time' points to  $t$ , and 'CPU, RAM, GPU power draw' points to the sum of power draws in the numerator.

# Measuring emissions

- First, measure power consumption:

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

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  $p_t$  result.

- Next, measure emissions:

# Measuring emissions

- First, measure power consumption:

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

Annotations: "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; "watts" points to the  $p_t$  result.

- Next, measure emissions:

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

# Measuring emissions

- First, measure power consumption:

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

Labels in the diagram:  
- PUE points to the PUE term in the numerator.  
- Running Time points to the  $t$  term in the numerator.  
- CPU, RAM, GPU power draw points to the  $(p_c + p_r + p_g)$  term in the numerator.  
- watts points to the  $p_t$  term on the left side of the equation.

- Next, measure emissions:

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

# Measuring emissions

- First, measure power consumption:

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

watts

- Next, measure emissions:

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

Power consumption of experiments

# Measuring emissions

- First, measure power consumption:

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

watts

- Next, measure emissions:

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

avg. CO<sub>2</sub>e (kg) per kWh where experiments took place

Power consumption of experiments

# Measuring emissions

- First, measure power consumption:

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

watts

- Next, measure emissions:

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

emissions

avg. CO<sub>2</sub>e (kg) per kWh where experiments took place

Power consumption of experiments

- Emissions of my search engine:

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

# Measuring emissions

- First, measure power consumption:

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

watts

- Next, measure emissions:

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

emissions

avg. CO<sub>2</sub>e (kg) per kWh where experiments took place

Power consumption of experiments

- Emissions of my search engine:

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

Power consumption of a single query

# Measuring emissions

- First, measure power consumption:

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

watts

- Next, measure emissions:

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

emissions

avg. CO<sub>2</sub>e (kg) per kWh where experiments took place

Power consumption of experiments

- Emissions of my search engine:

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

No. queries issued per unit time

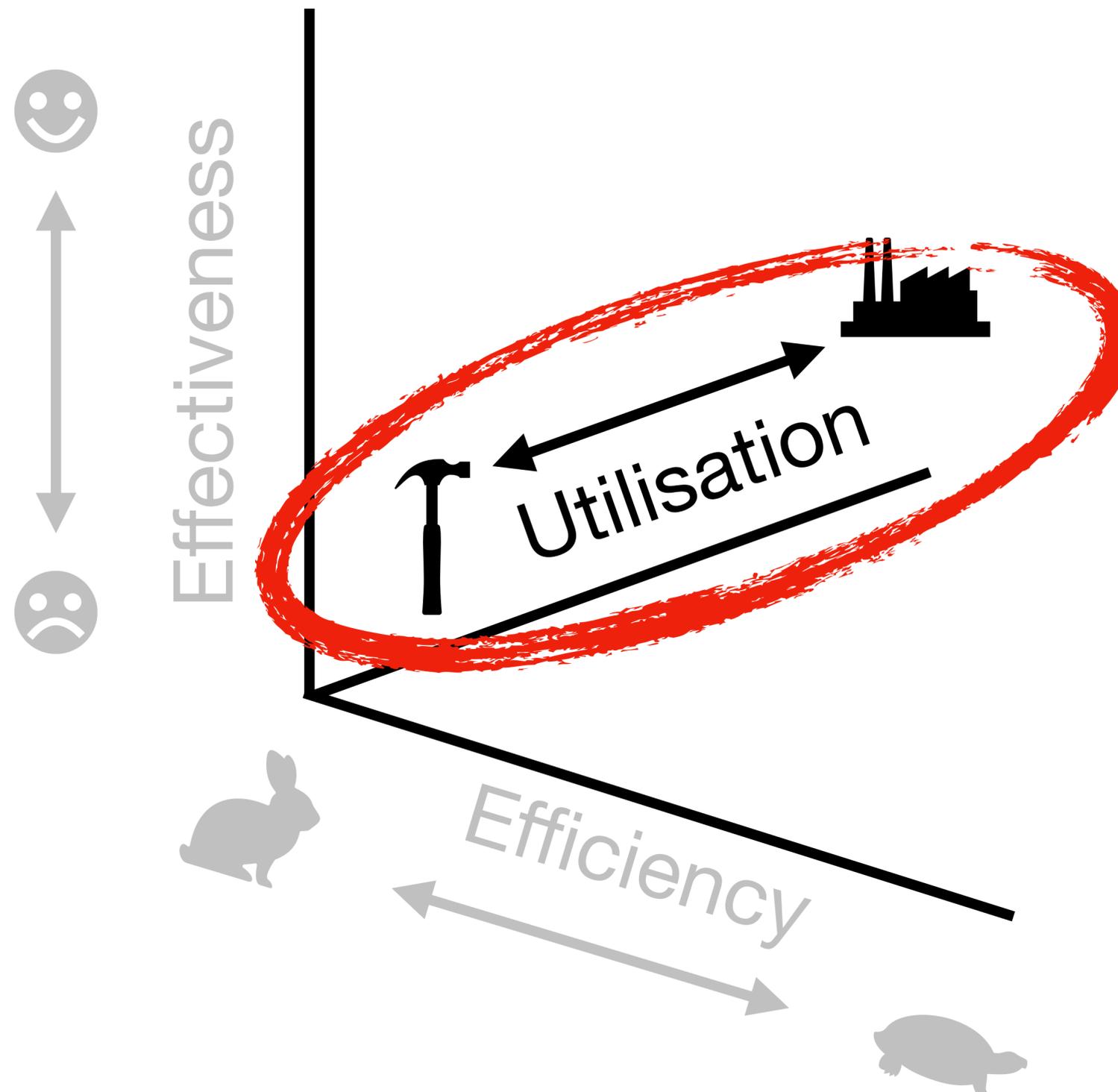
Power consumption of a single query

# Measuring power and emissions in practice

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

```
from codecarbon import EmissionsTracker

tracker = EmissionsTracker()
tracker.start()
# Experiment code goes here
tracker.stop()
```



~~Okay, so what does this mean for IR?~~

~~Okay, so how can I measure this?~~

**Okay, so show me what it means in IR research practice!**

# Experimental Setup Overview

- Methods:
  - BM25
  - LambdaMART
  - DPR
  - monoBERT
  - uniCOIL
  - TILDEv2

# Experimental Setup Overview

- Methods:

- BM25

- LambdaMART

Non-neural

- DPR

- monoBERT

- uniCOIL

- TILDEv2

# Experimental Setup Overview

- Methods:

- BM25
- LambdaMART

Non-neural

- DPR
- monoBERT
- uniCOIL
- TILDEv2

“Neural”

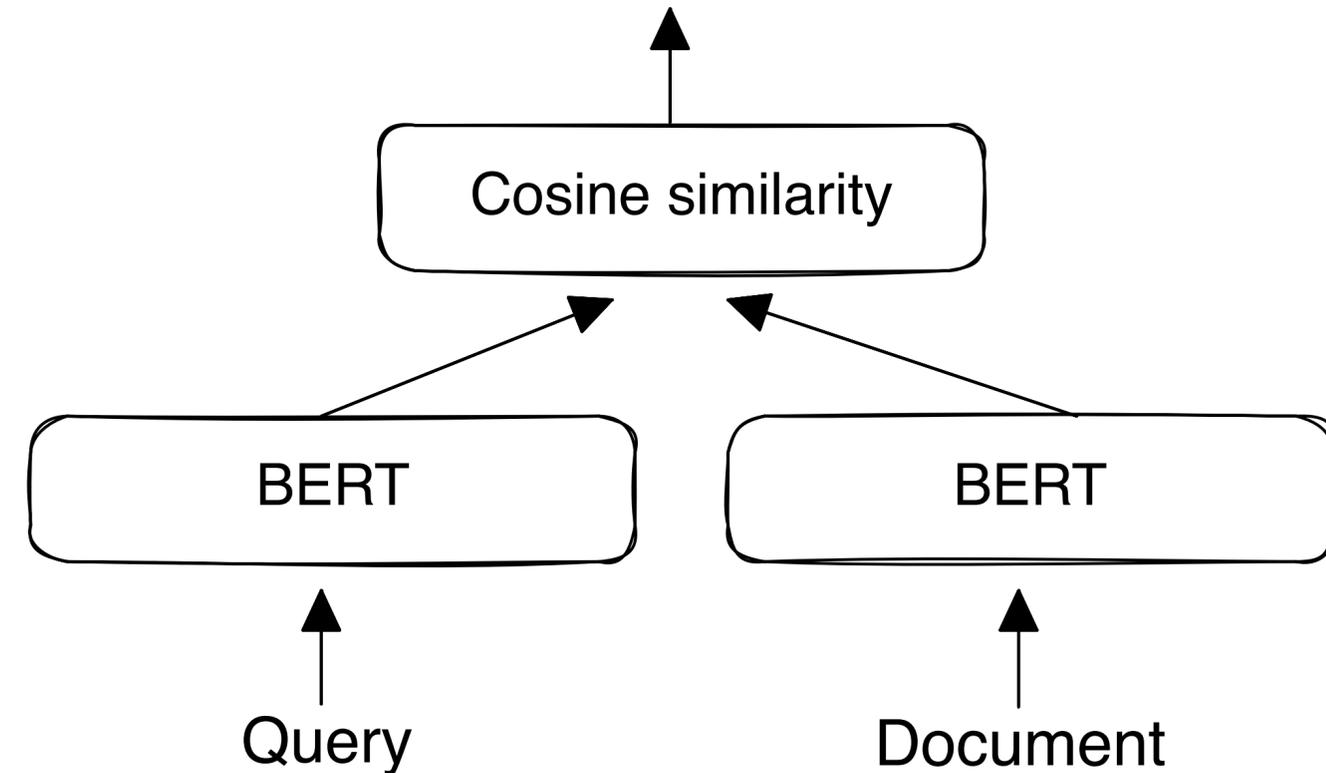
# Experimental Setup Overview

- Methods:

- BM25
- LambdaMART
- DPR
- monoBERT
- uniCOIL
- TILDEv2

Non-neural

Dense retriever (bi-encoder)



# Experimental Setup Overview

- Methods:

- BM25

- LambdaMART

- DPR

- monoBERT

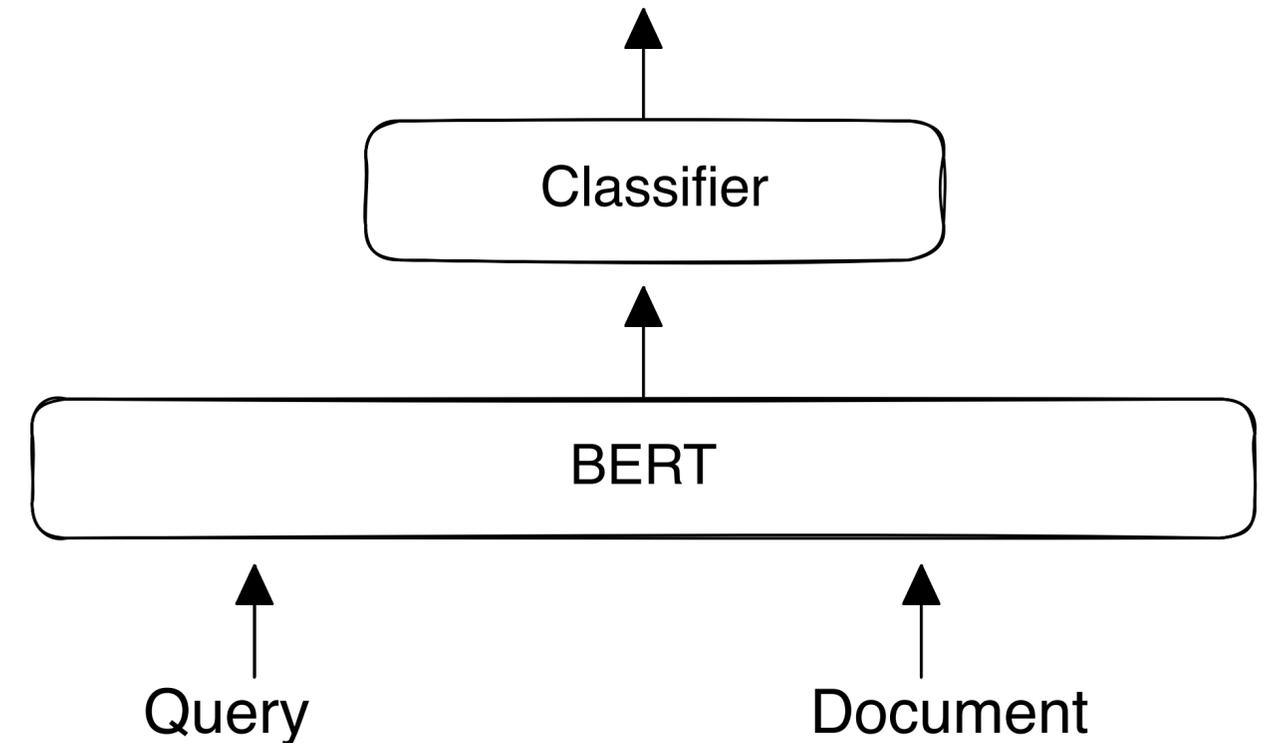
- uniCOIL

- TILDEv2

Non-neural

Dense retriever (bi-encoder)

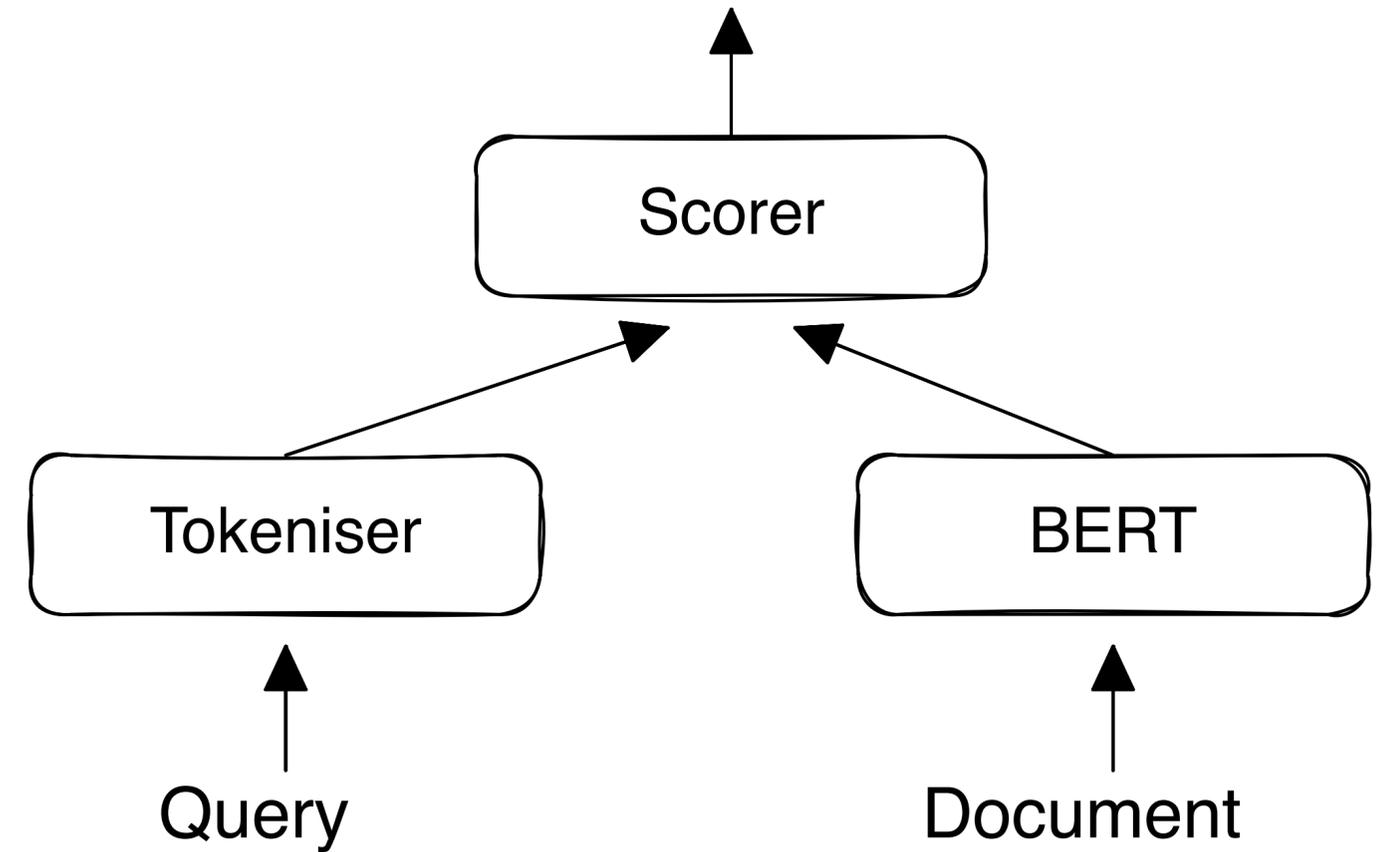
BERT (cross-encoder)



# Experimental Setup Overview

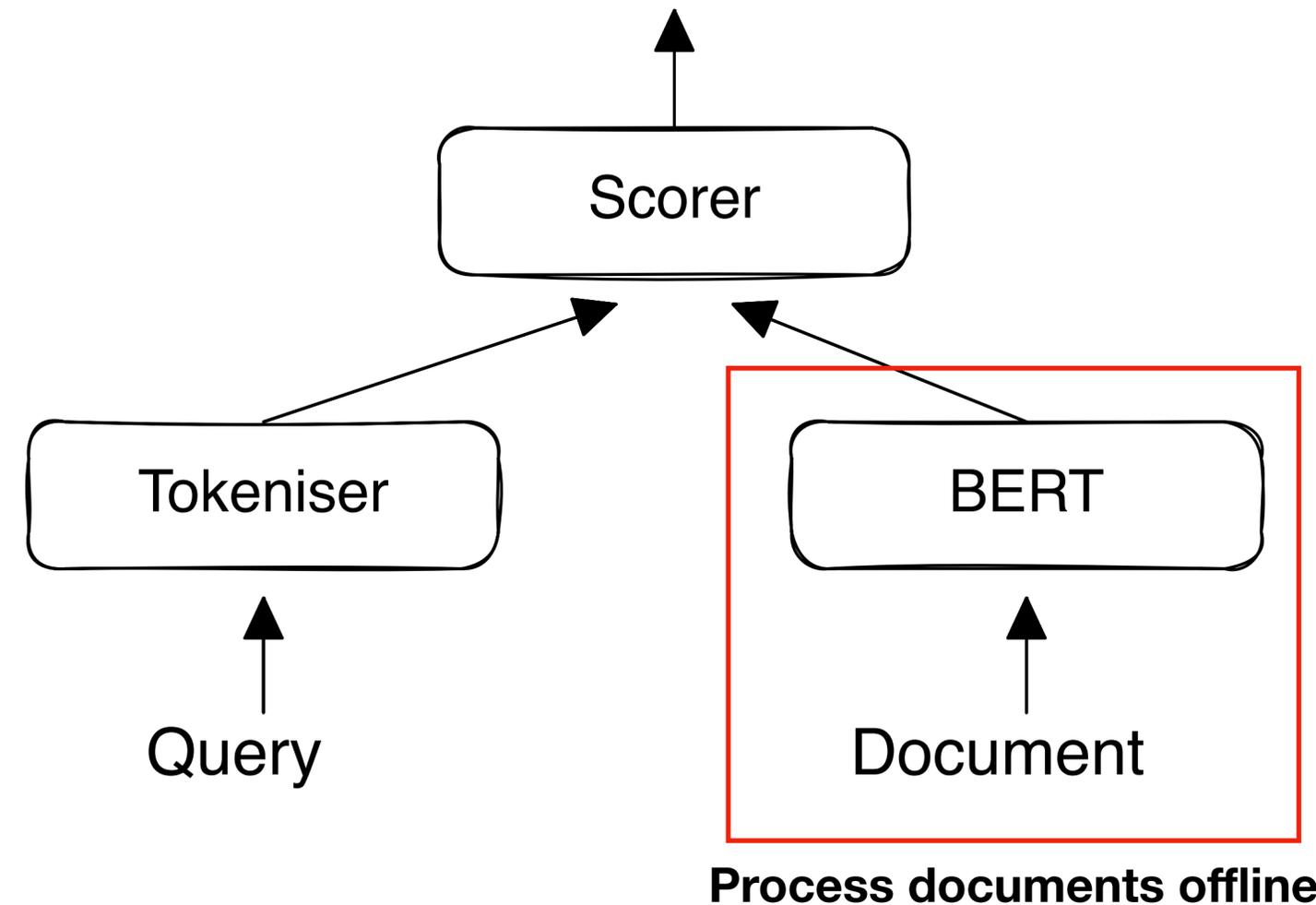
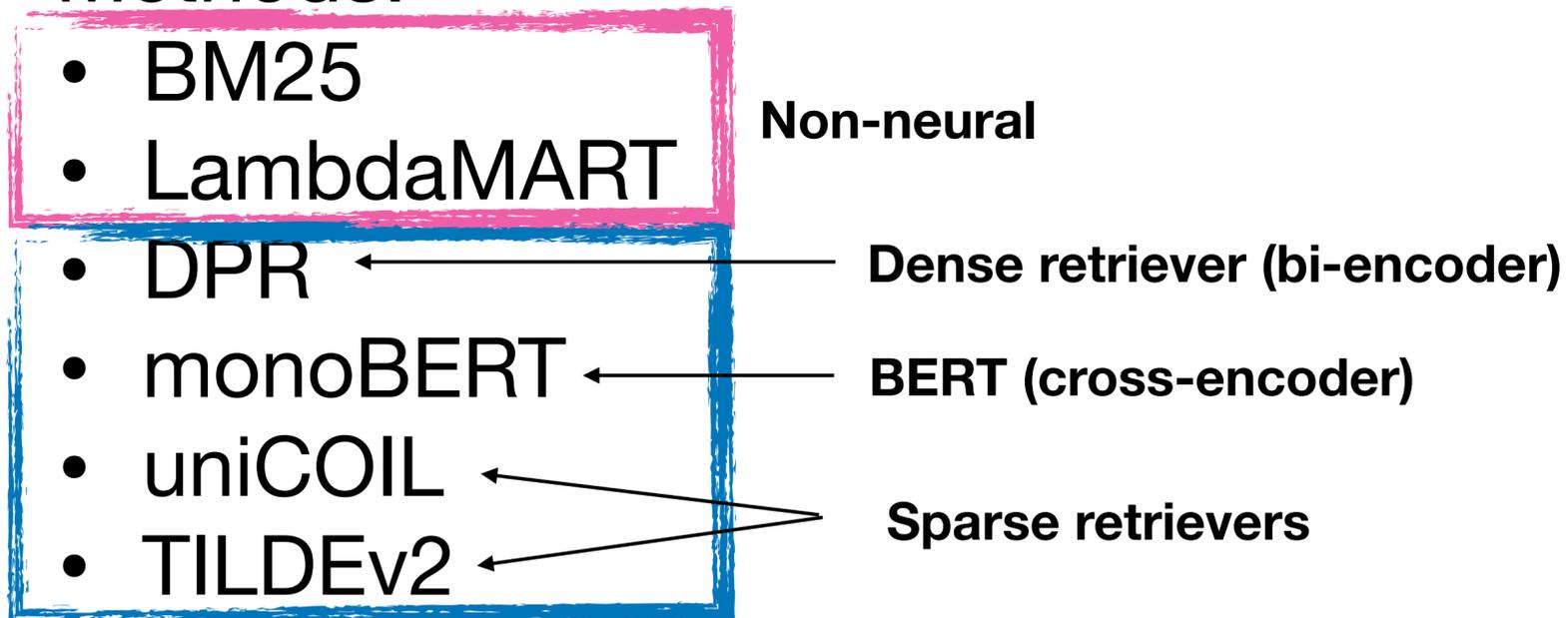
- Methods:

- BM25
- LambdaMART
- DPR ← Dense retriever (bi-encoder)
- monoBERT ← BERT (cross-encoder)
- uniCOIL ← Sparse retrievers
- TILDEv2 ← Sparse retrievers



# Experimental Setup Overview

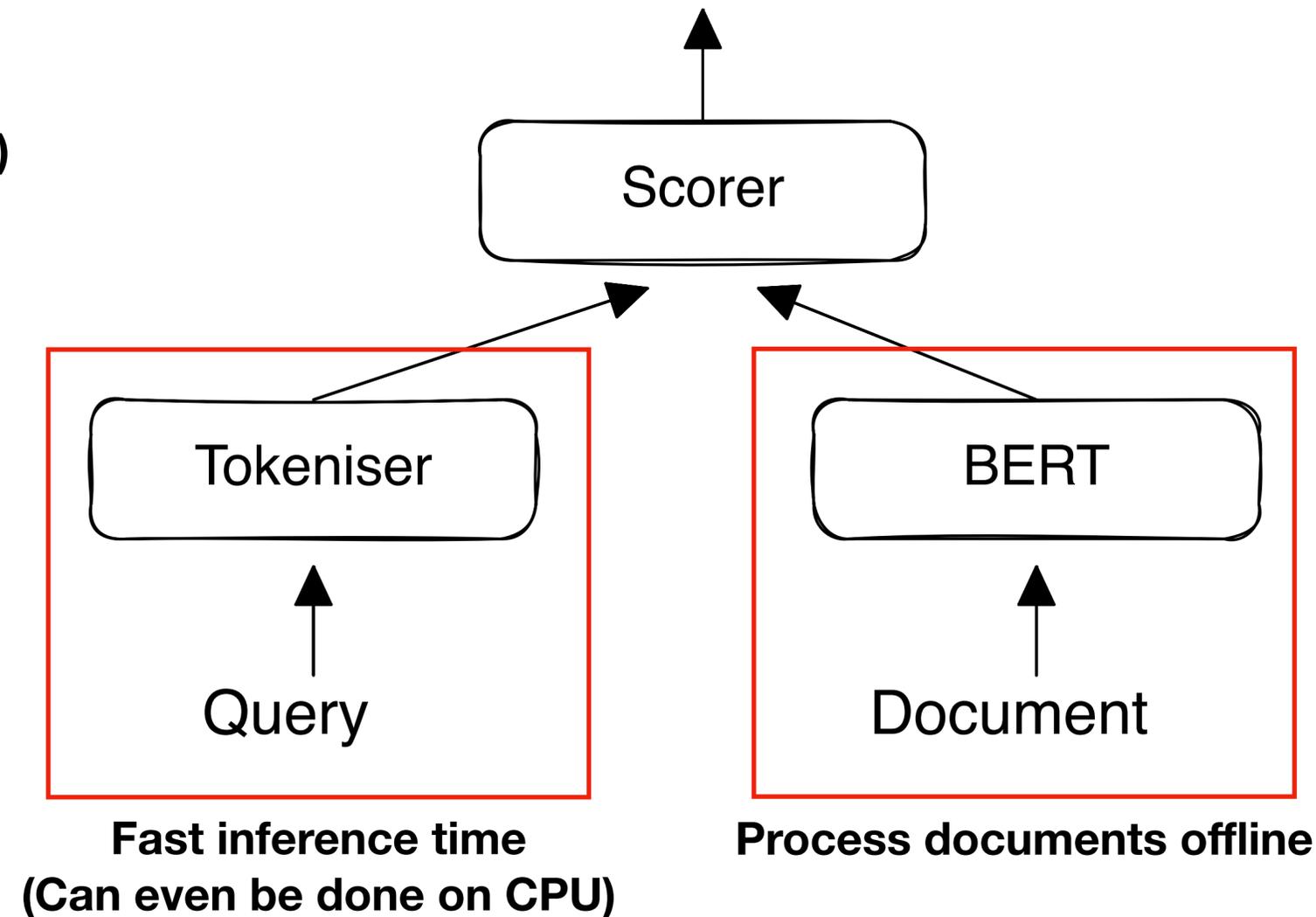
- Methods:



# Experimental Setup Overview

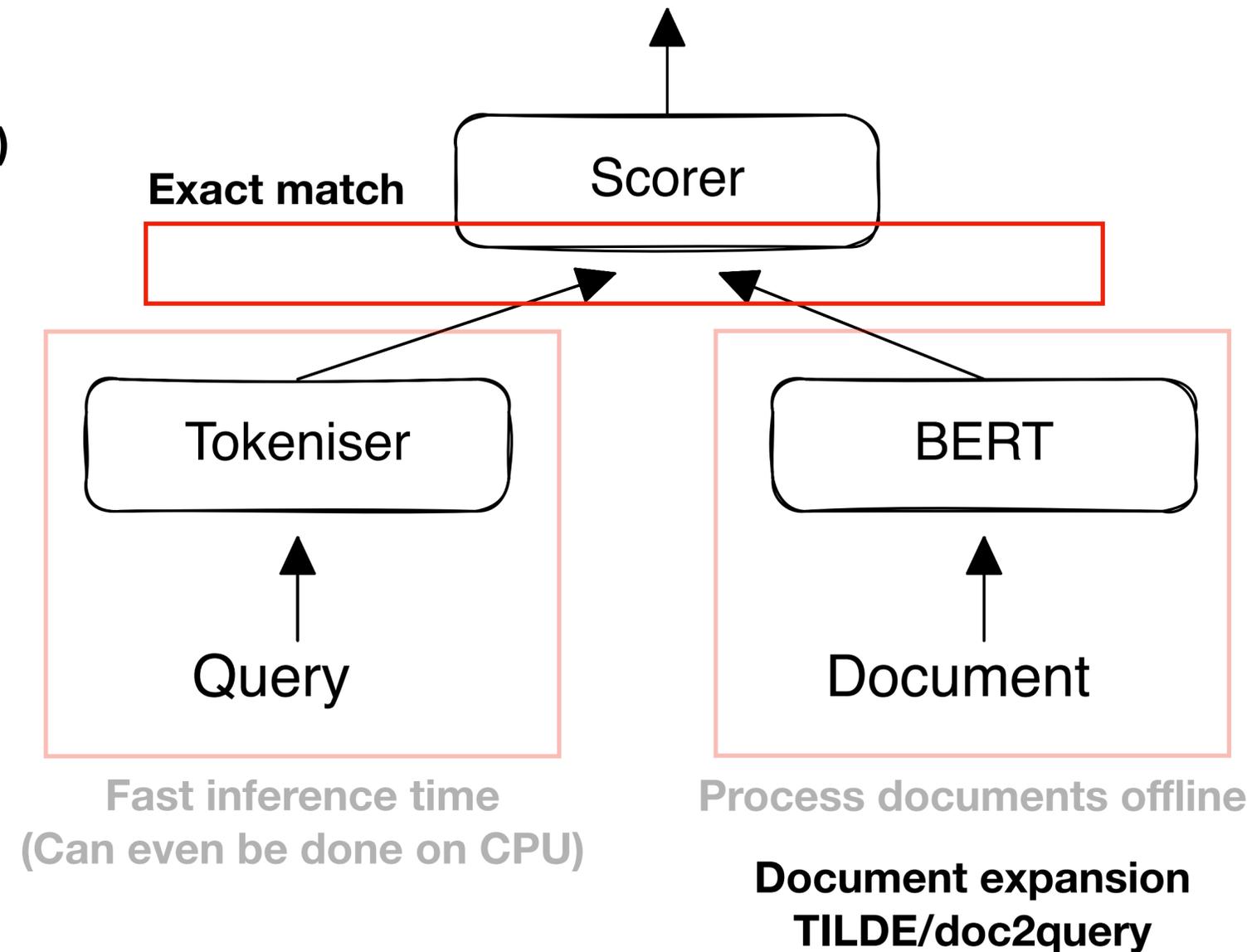
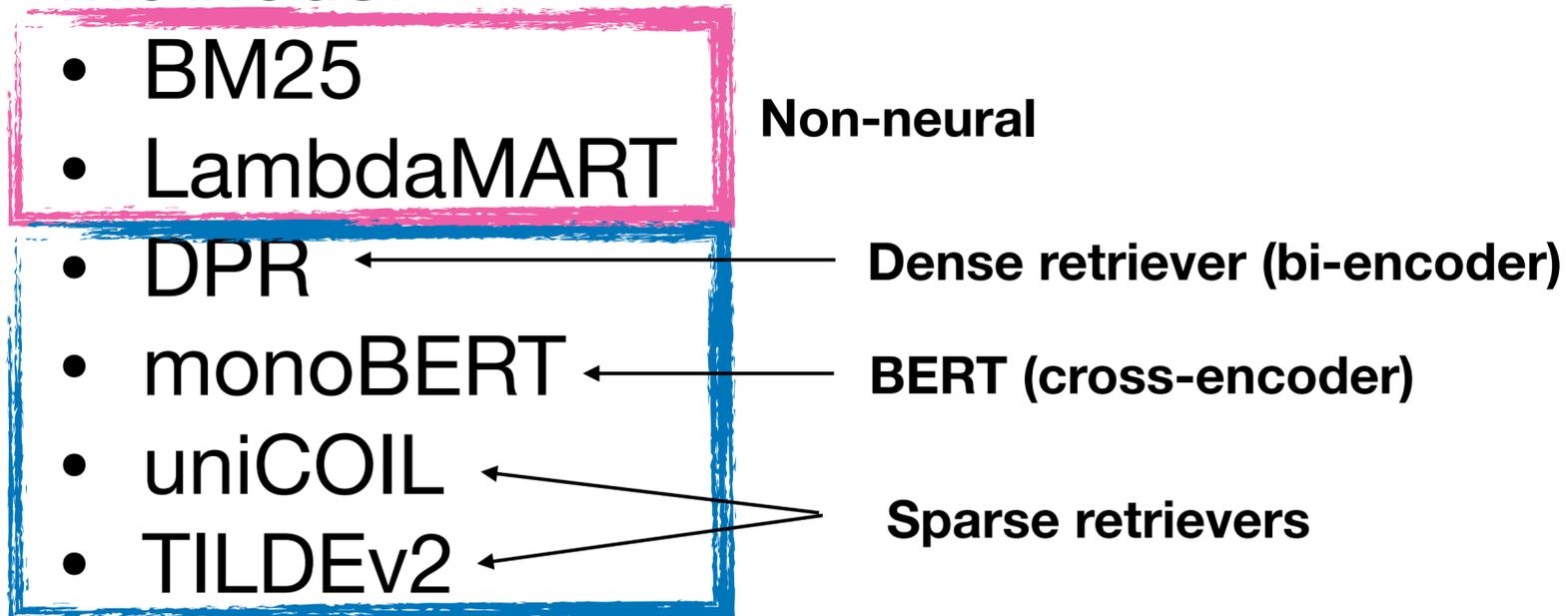
- Methods:

- BM25
- LambdaMART
- DPR ← Dense retriever (bi-encoder)
- monoBERT ← BERT (cross-encoder)
- uniCOIL ← Sparse retrievers
- TILDEv2 ← Sparse retrievers



# Experimental Setup Overview

- Methods:



# Experimental Setup Overview

- Methods:

- BM25
- LambdaMART

Non-neural

- DPR
- monoBERT
- uniCOIL
- TILDEv2

“Neural”

- Collection:

- MSMARCOv1

- Experiments:

# Experimental Setup Overview

- Methods:

- BM25
- LambdaMART

Non-neural

- DPR
- monoBERT
- uniCOIL
- TILDEv2

“Neural”

- Collection:

- MSMARCOv1

- Experiments:

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

# Experimental Setup Overview

- Methods:

- BM25
- LambdaMART

Non-neural

- DPR
- monoBERT
- uniCOIL
- TILDEv2

“Neural”

- Collection:

- MSMARCOv1

- Experiments:

- How many emissions do these methods produce to obtain an experimental result?
- What are the effectiveness-utilisation trade-offs of these methods?

# Experimental Setup Overview

- Methods:

- BM25
- LambdaMART

Non-neural

- DPR
- monoBERT
- uniCOIL
- TILDEv2

“Neural”

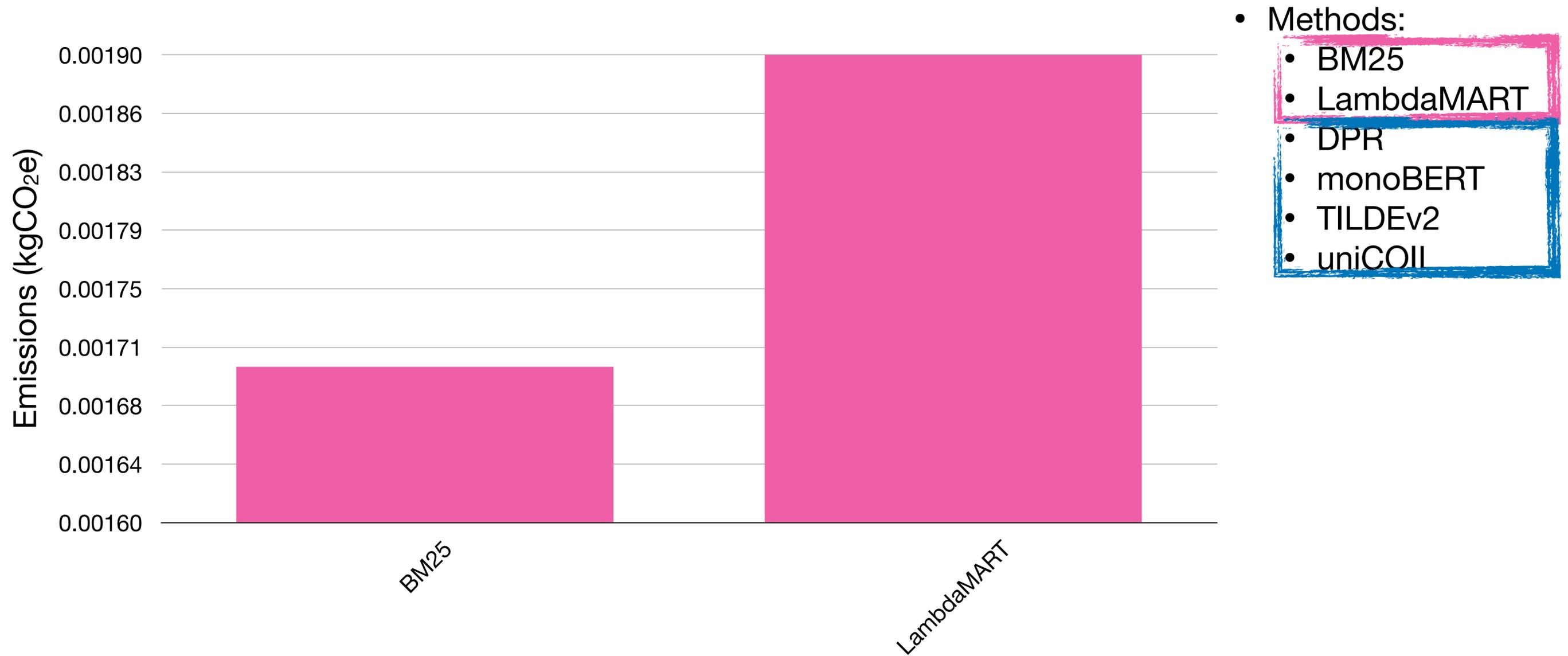
- Collection:

- MSMARCOv1

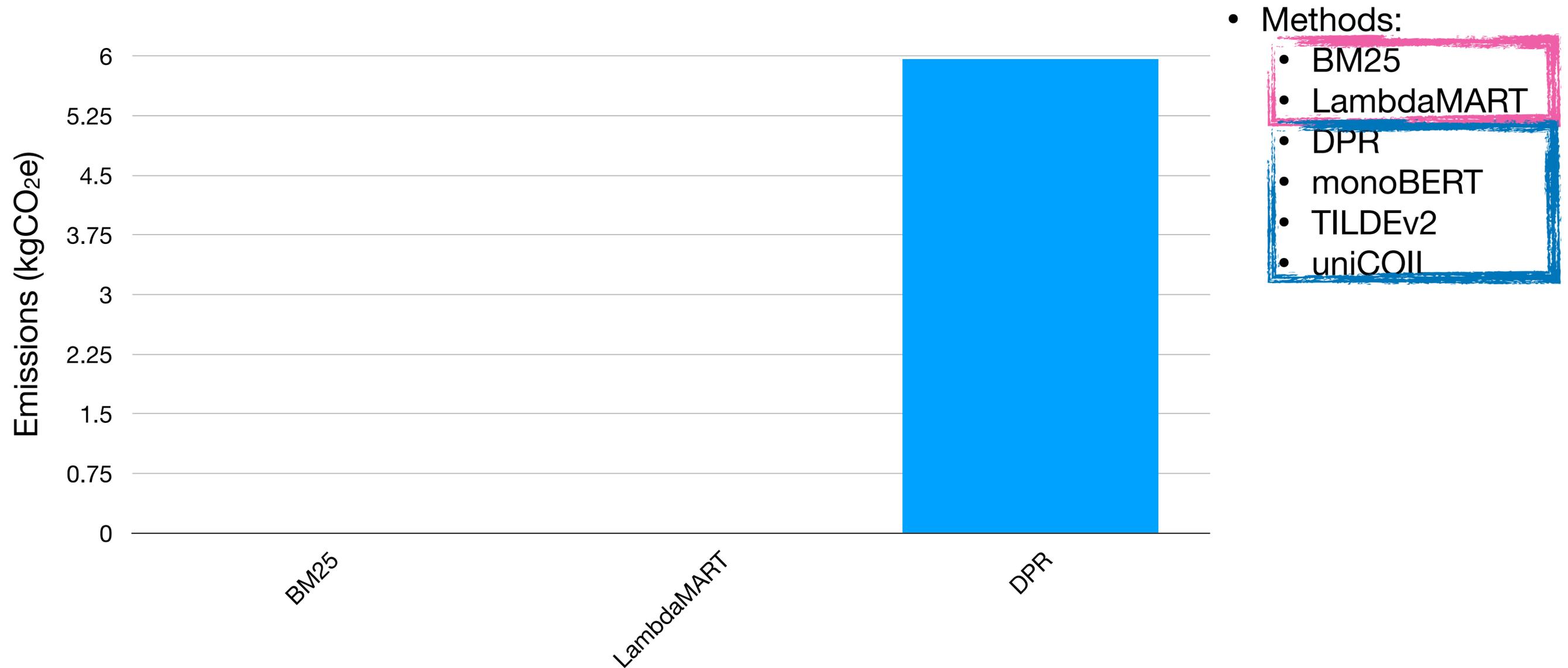
- Experiments:

- How many emissions do these methods produce to obtain an experimental result?
- What are the effectiveness-utilisation trade-offs of these methods?

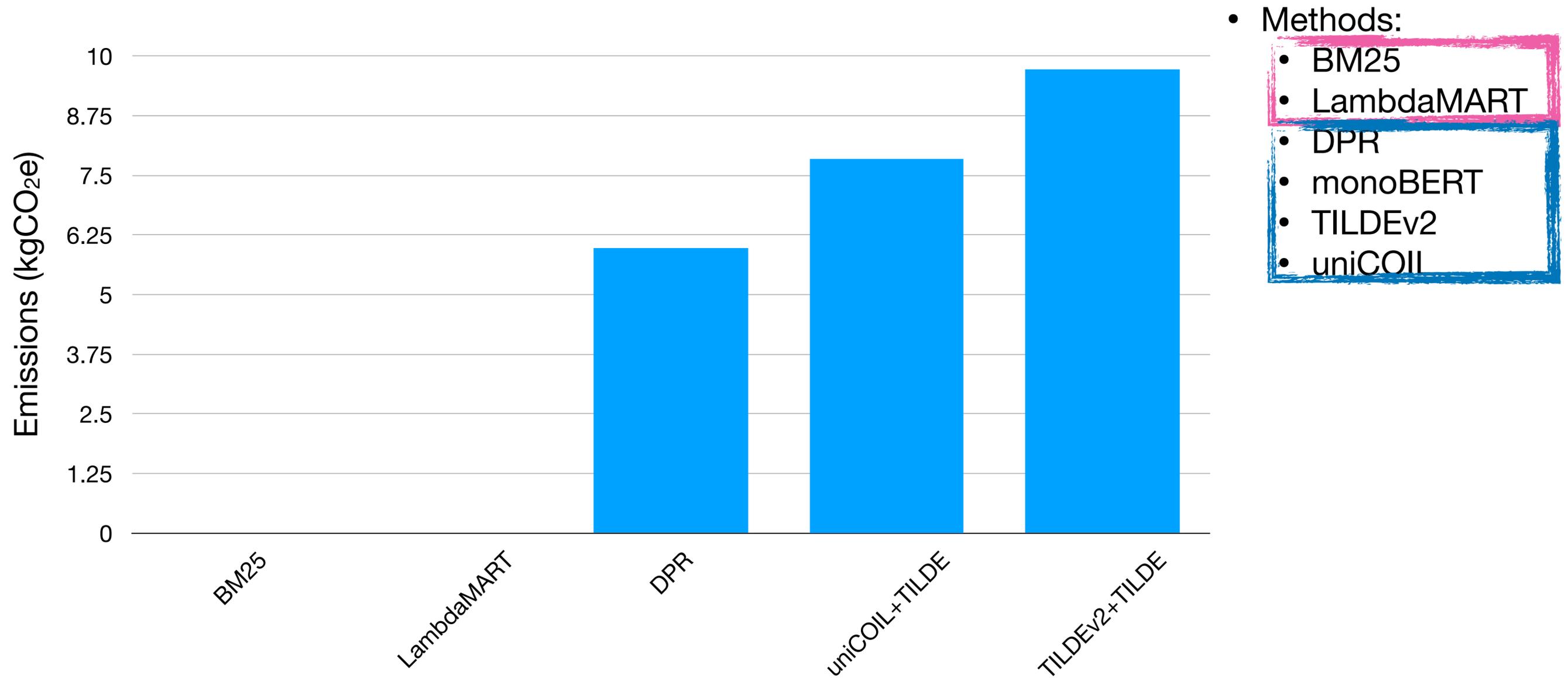
# How many emissions do these methods produce to obtain an experimental result?



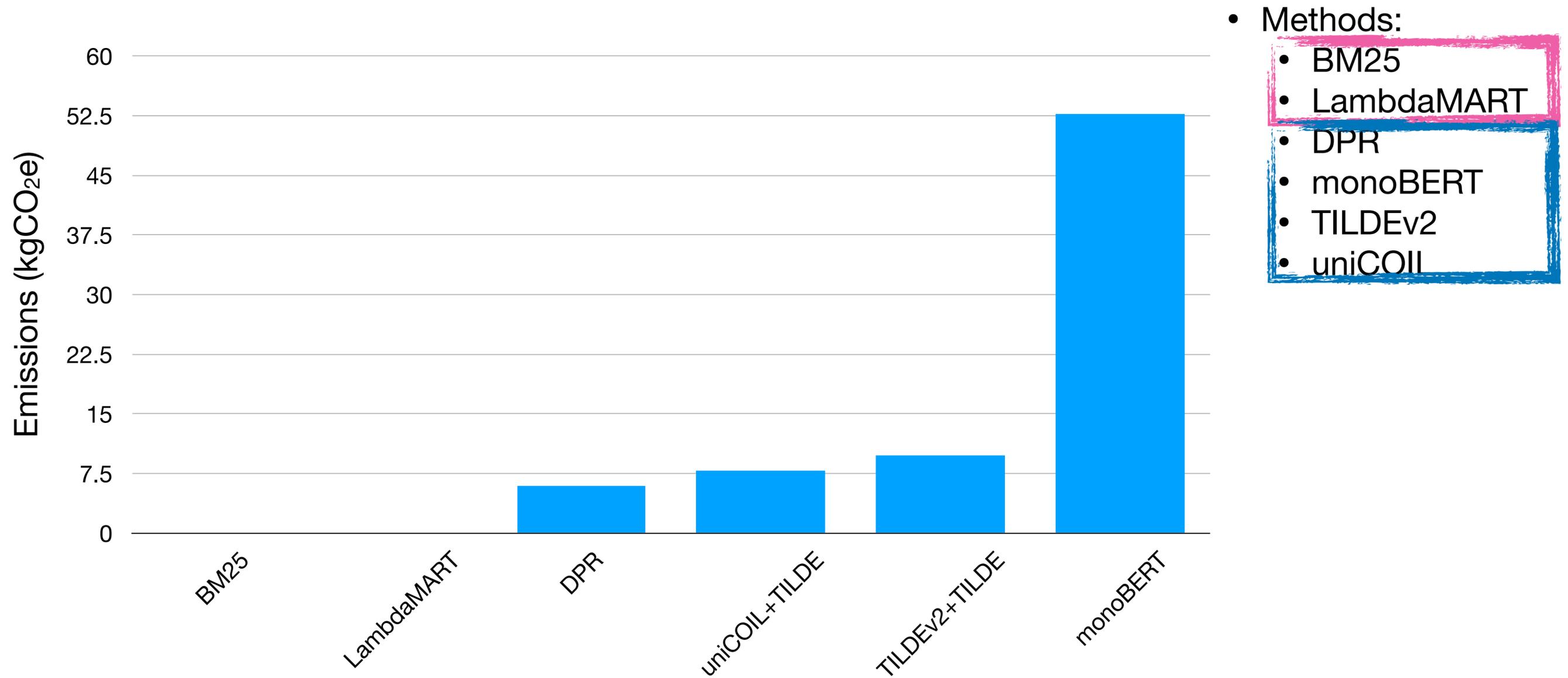
# How many emissions do these methods produce to obtain an experimental result?



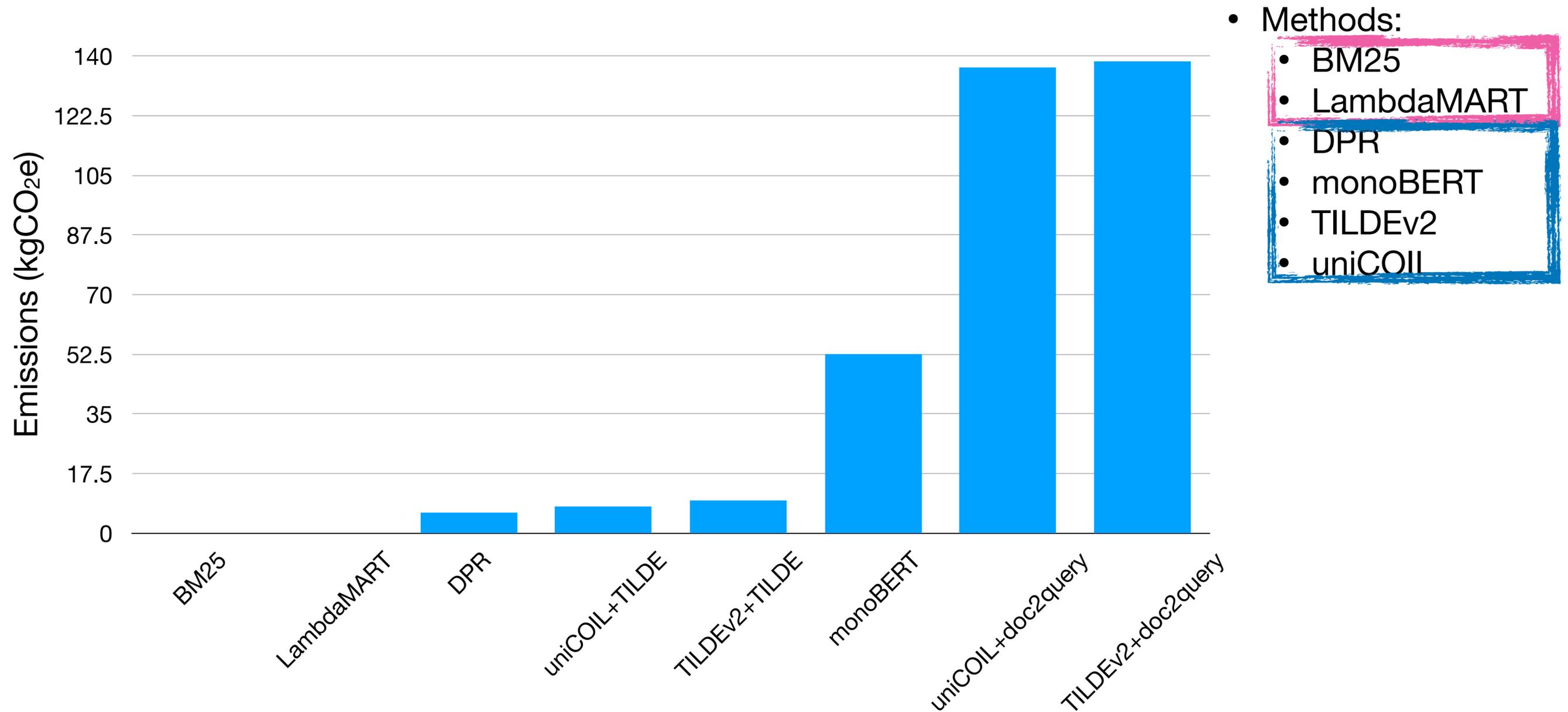
# How many emissions do these methods produce to obtain an experimental result?



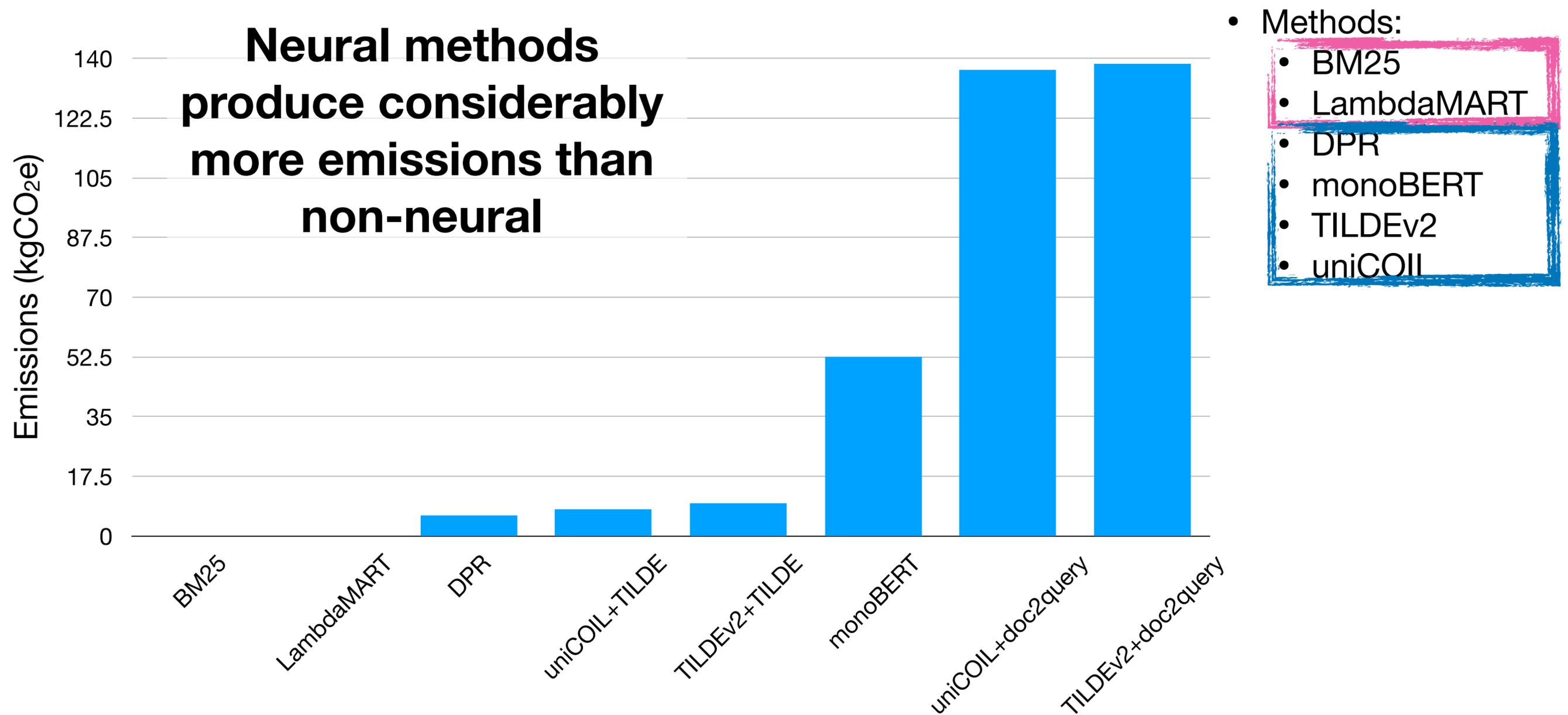
# How many emissions do these methods produce to obtain an experimental result?



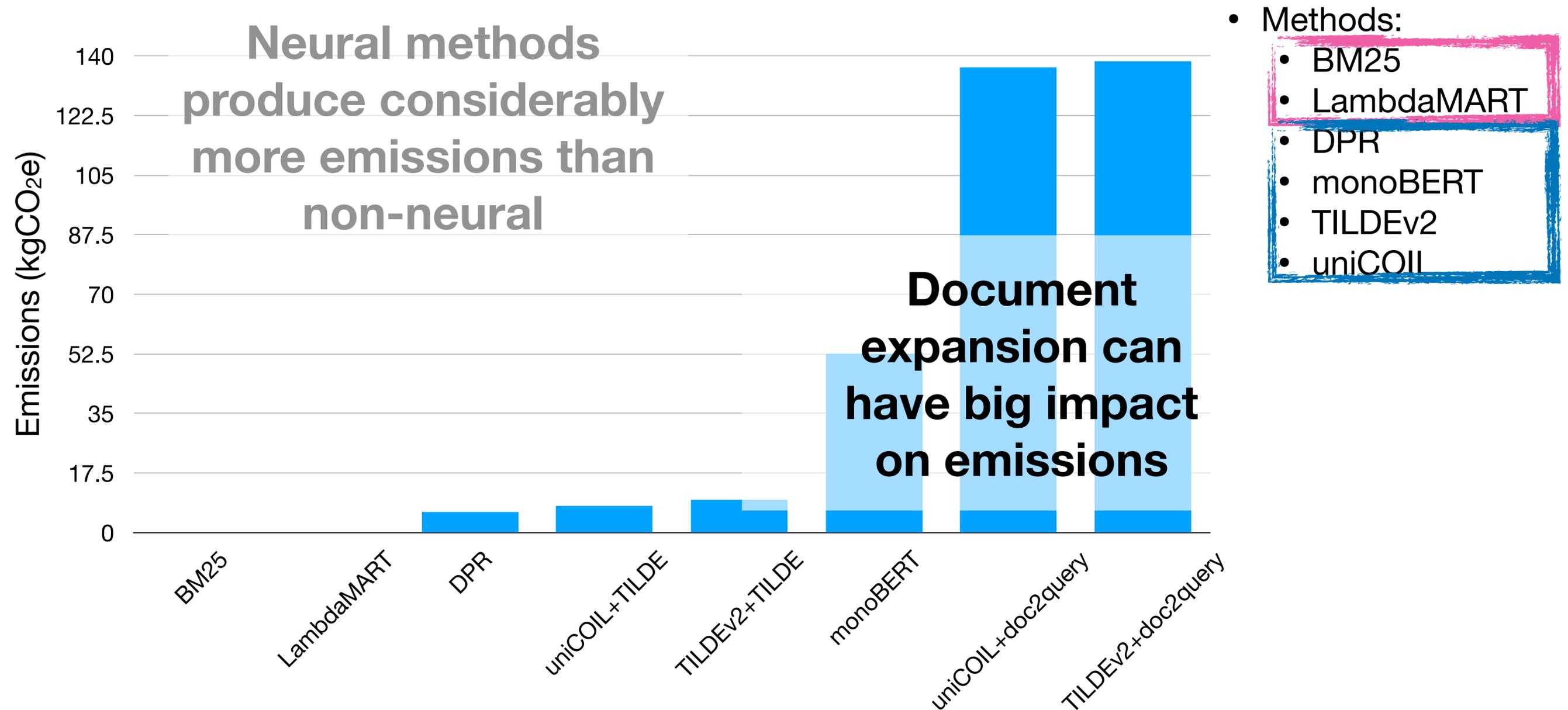
# How many emissions do these methods produce to obtain an experimental result?



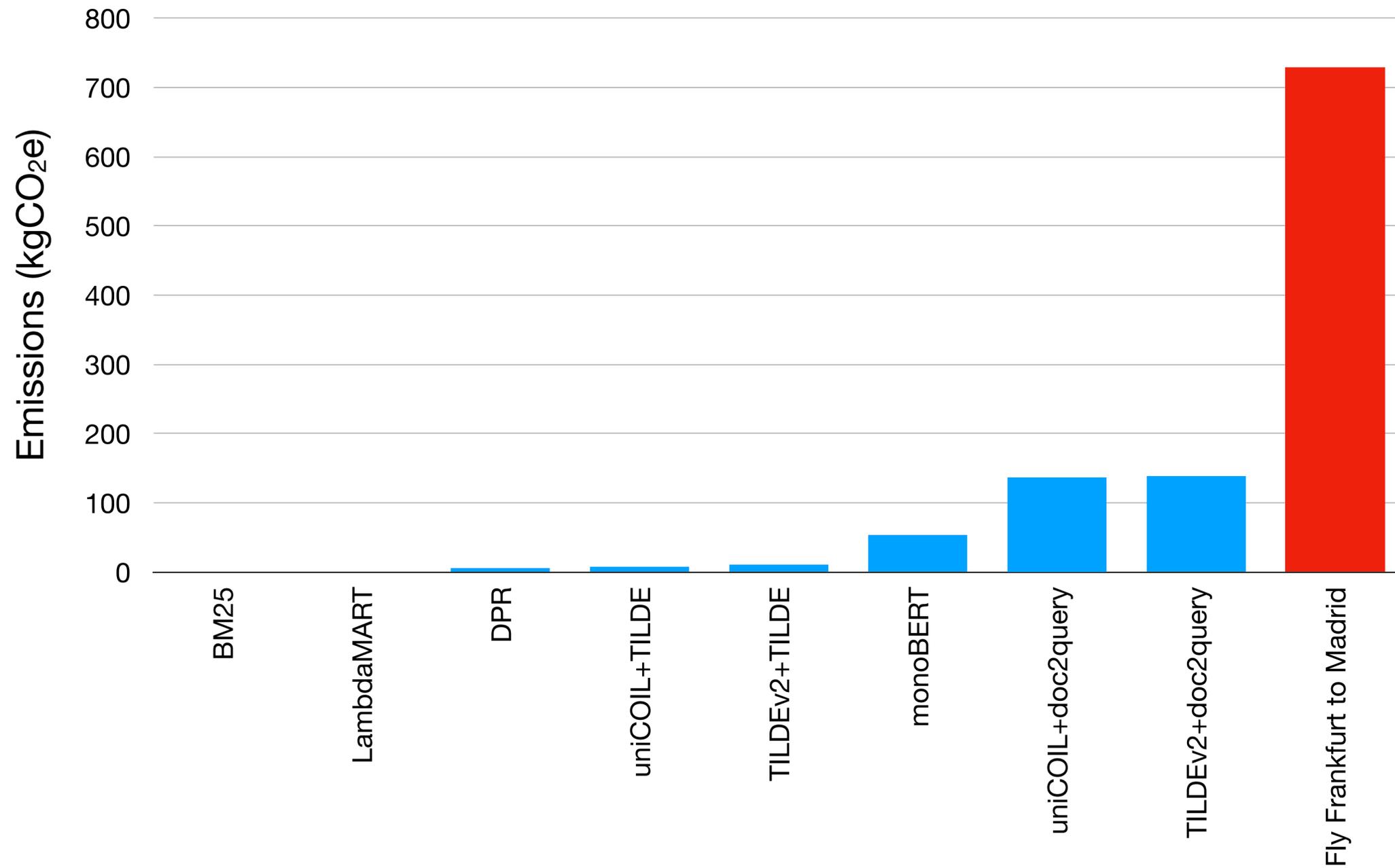
# How many emissions do these methods produce to obtain an experimental result?



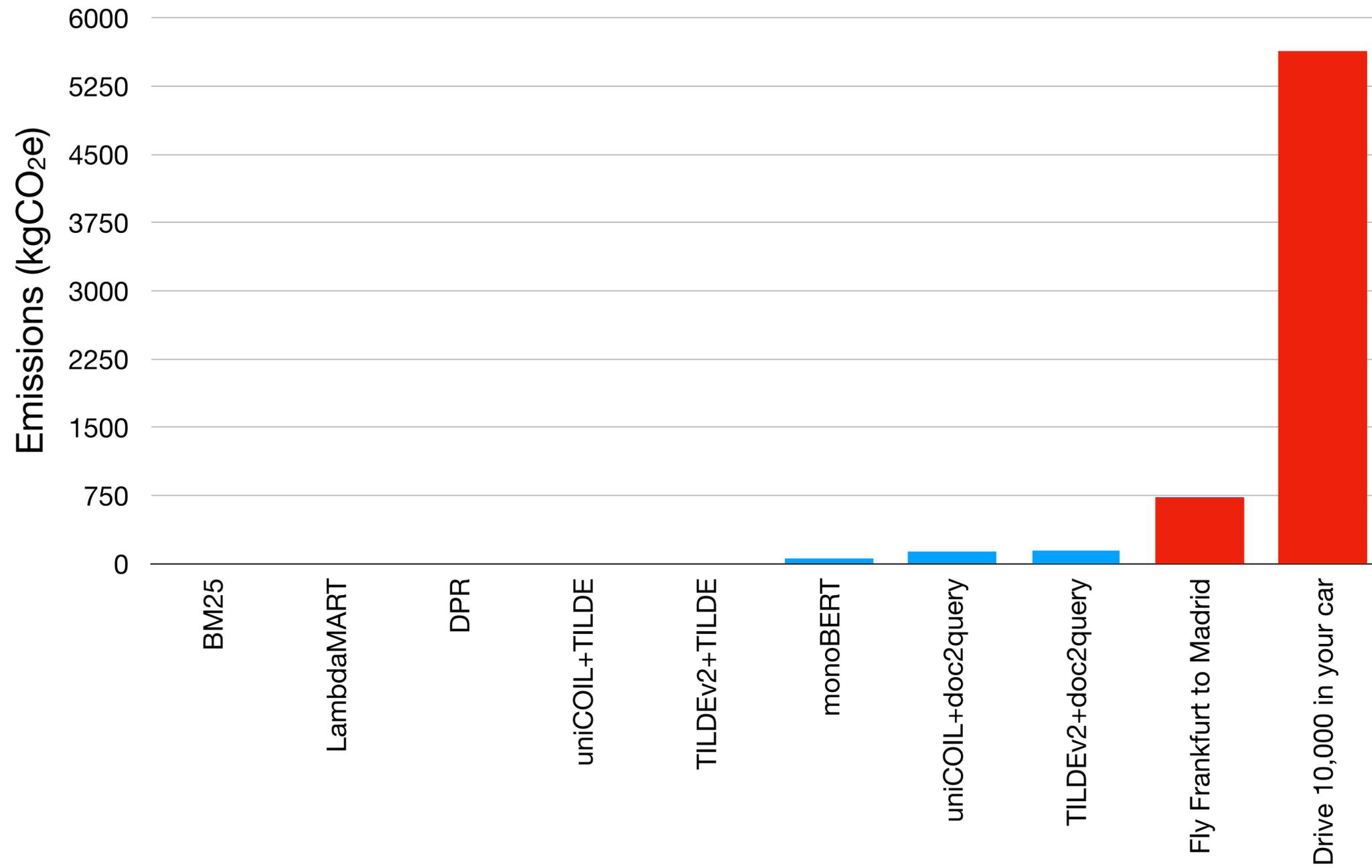
# How many emissions do these methods produce to obtain an experimental result?



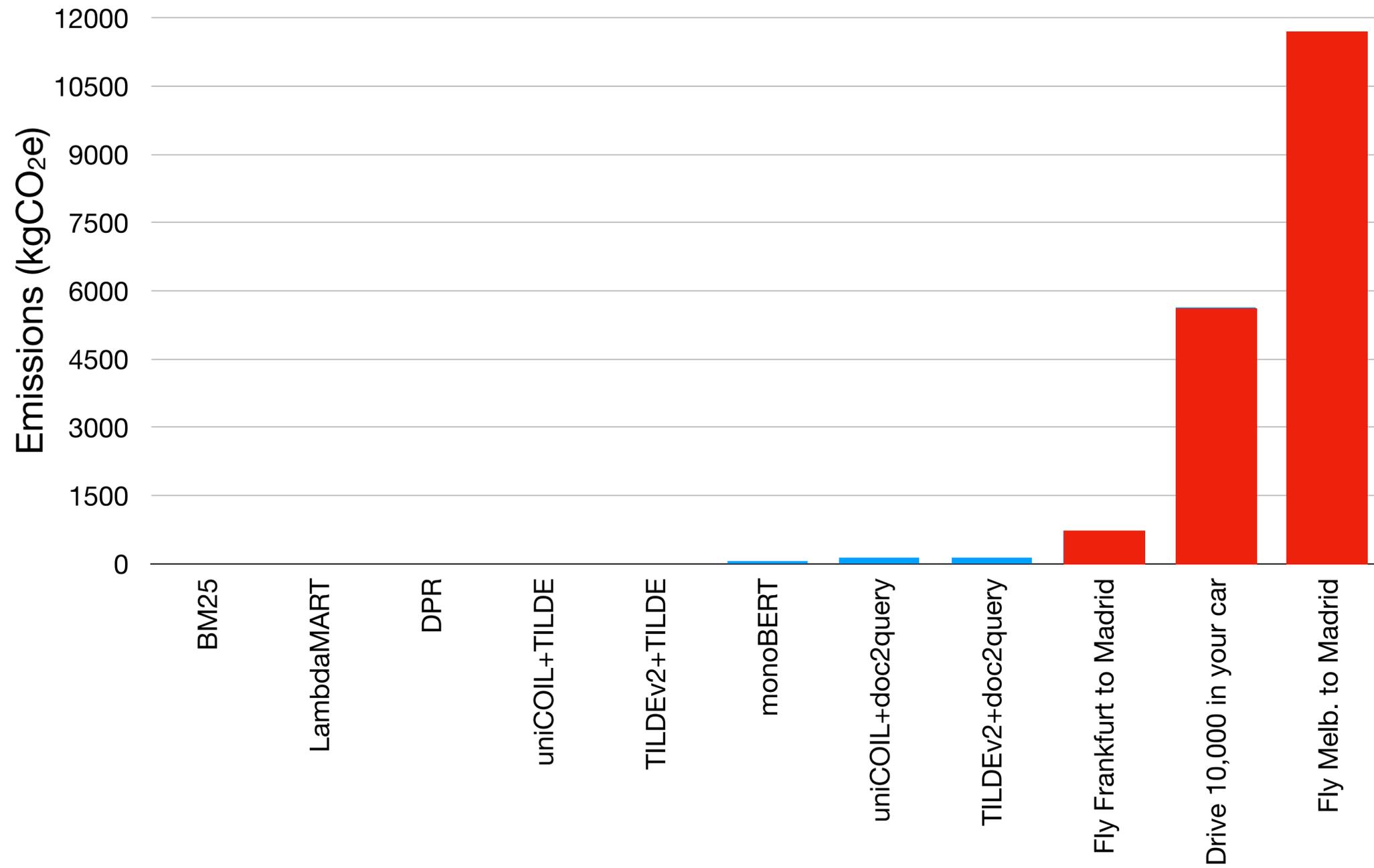
# How many emissions do these methods produce to obtain an experimental result?



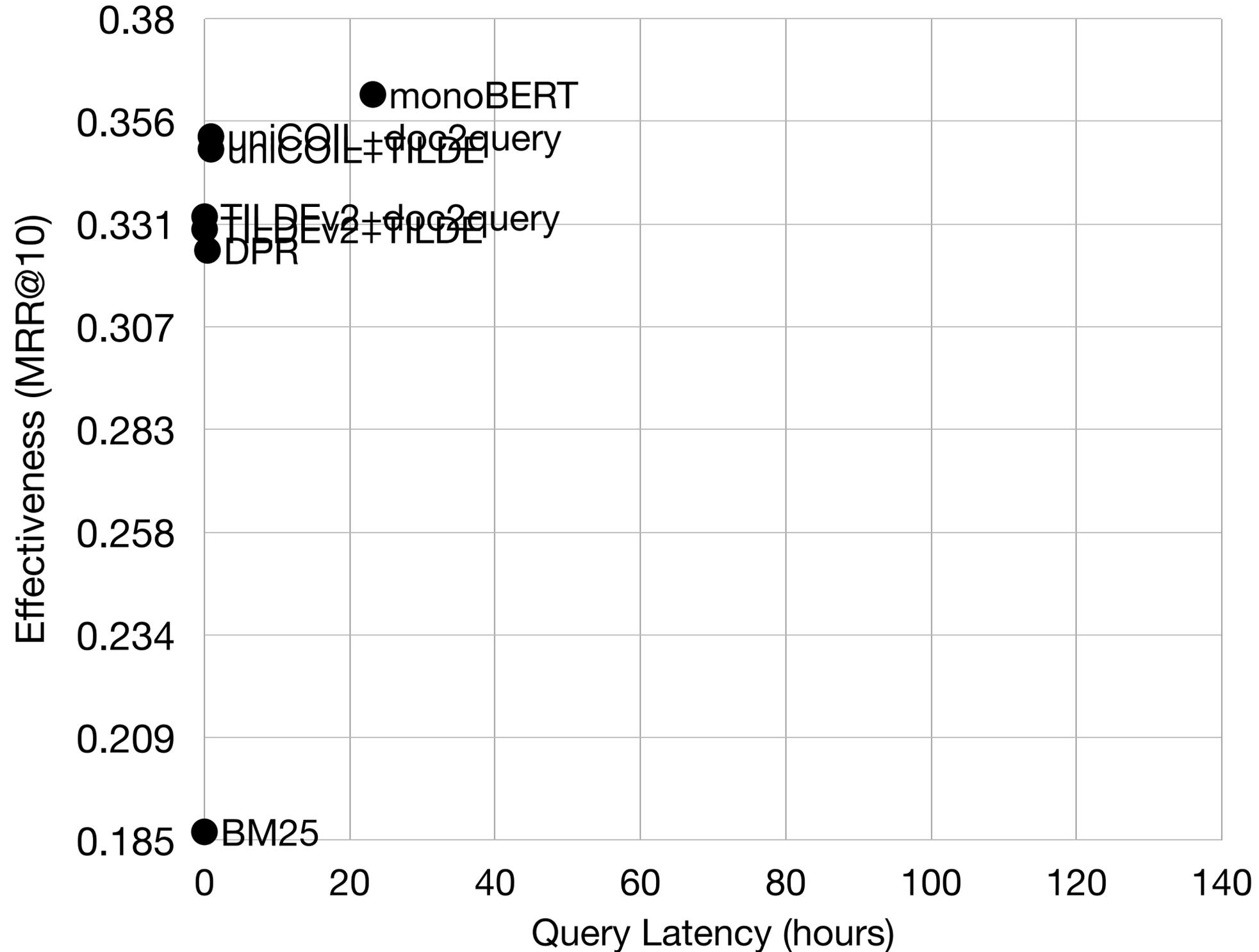
# How many emissions do these methods produce to obtain an experimental result?



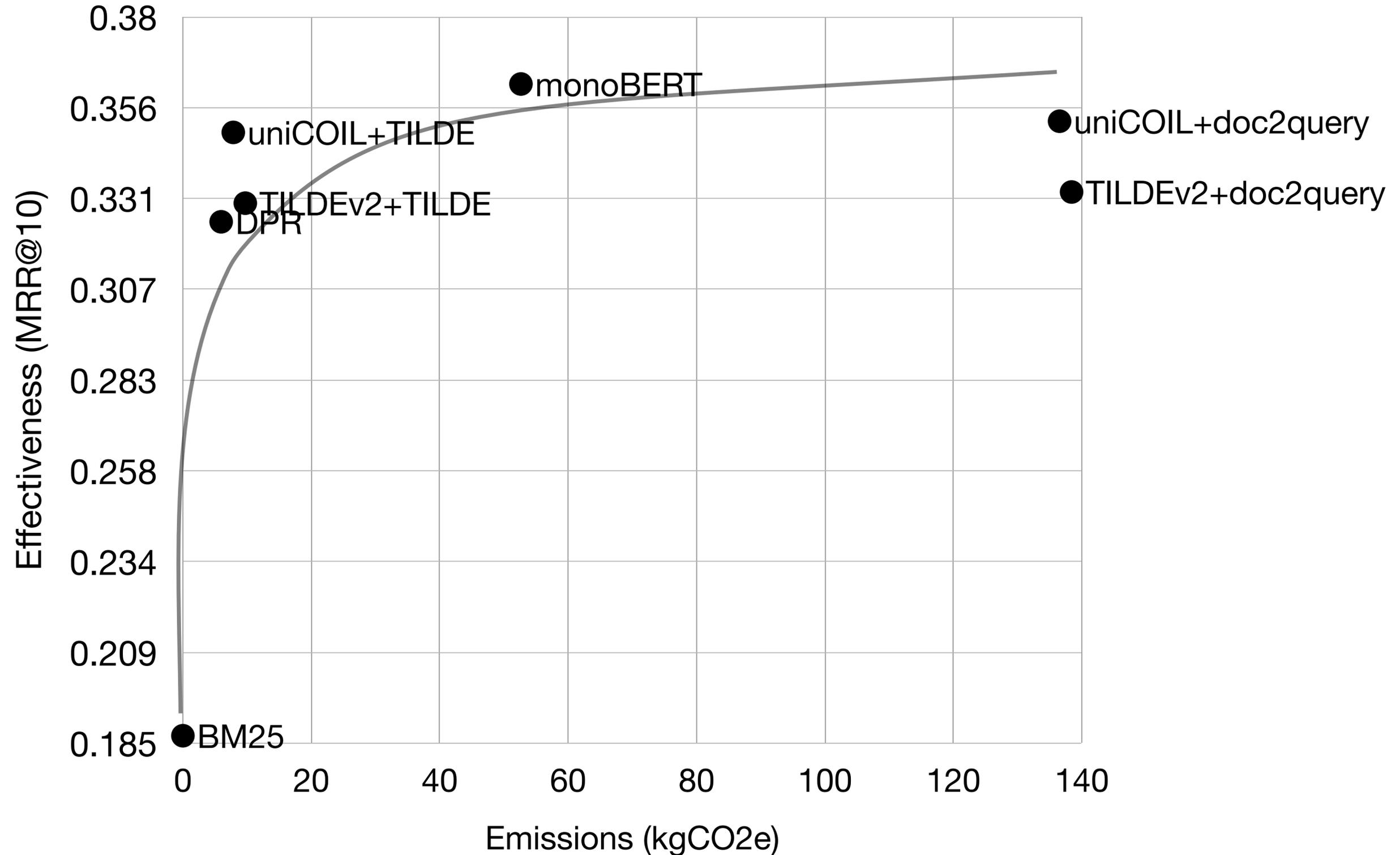
# How many emissions do these methods produce to obtain an experimental result?



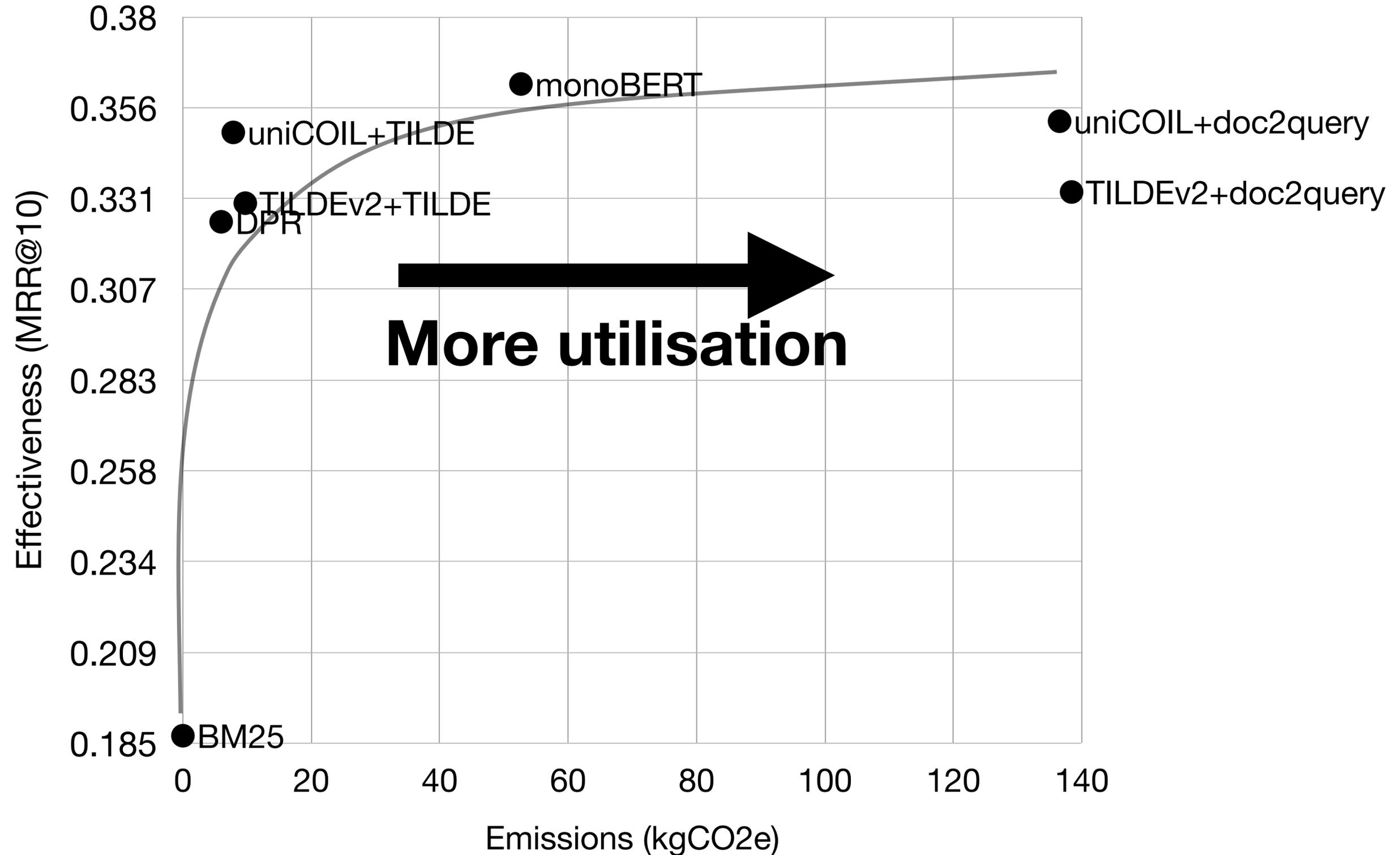
# What are the effectiveness-utilisation trade-offs of these methods?



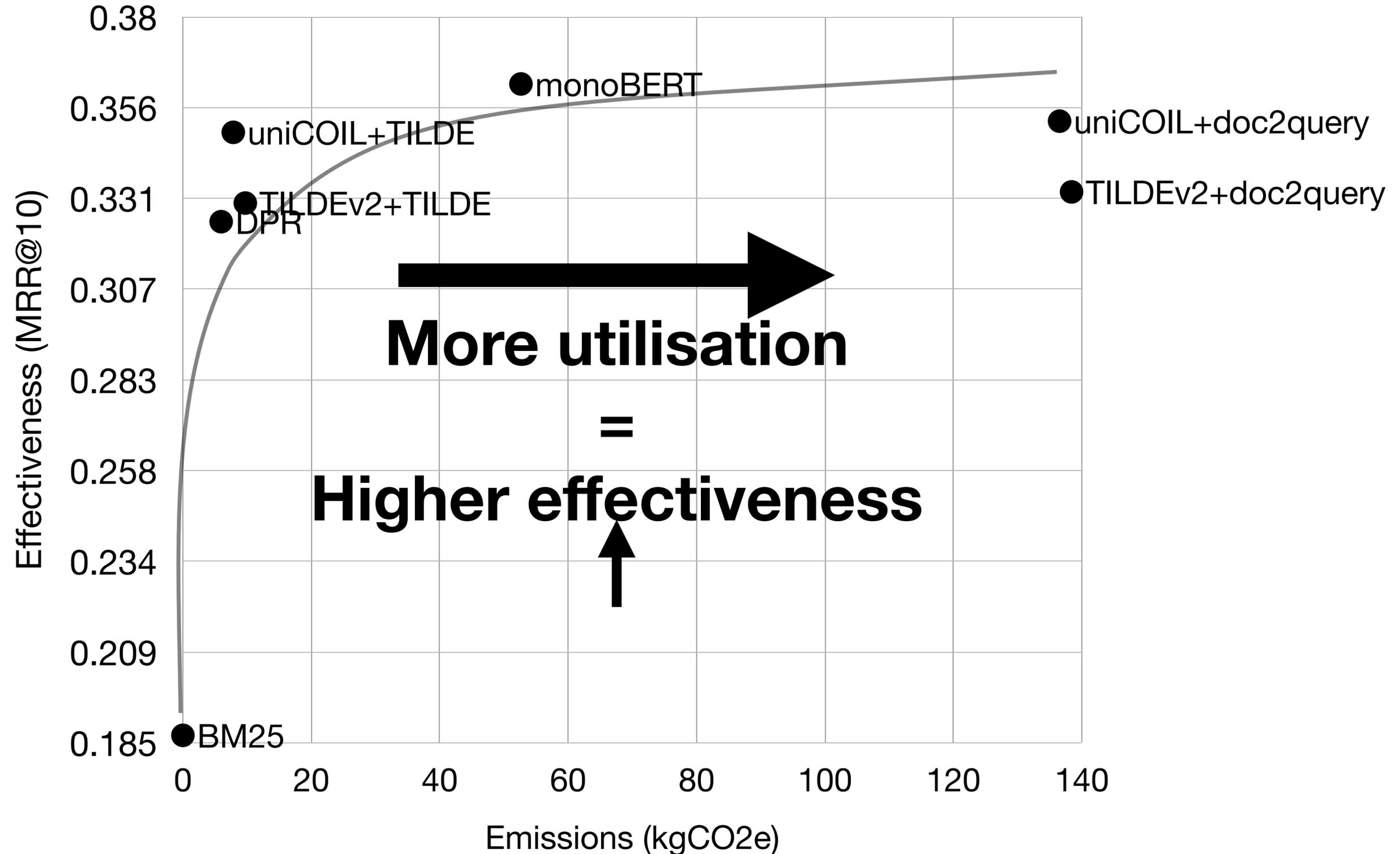
# What are the effectiveness-utilisation trade-offs of these methods?



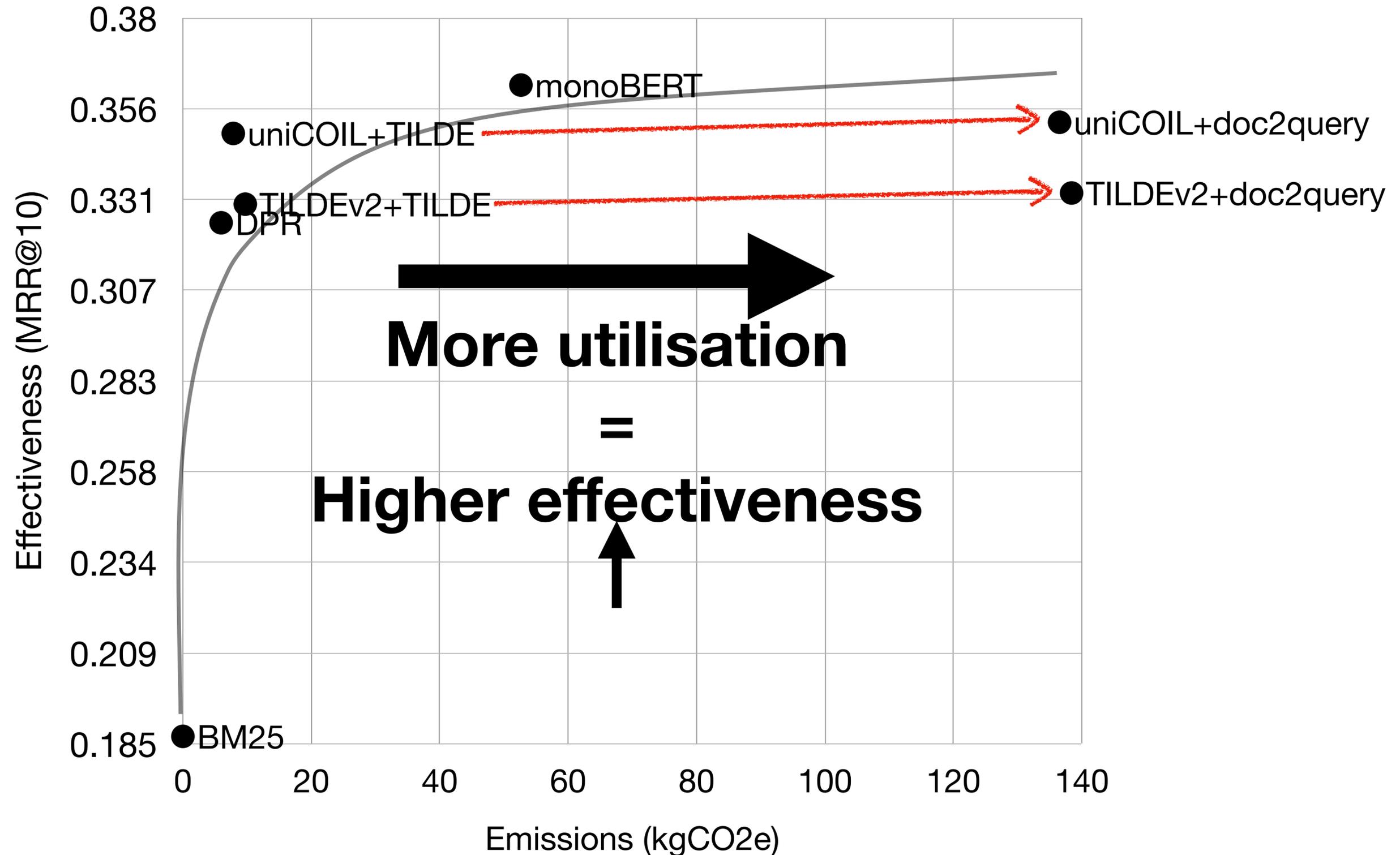
# What are the effectiveness-utilisation trade-offs of these methods?



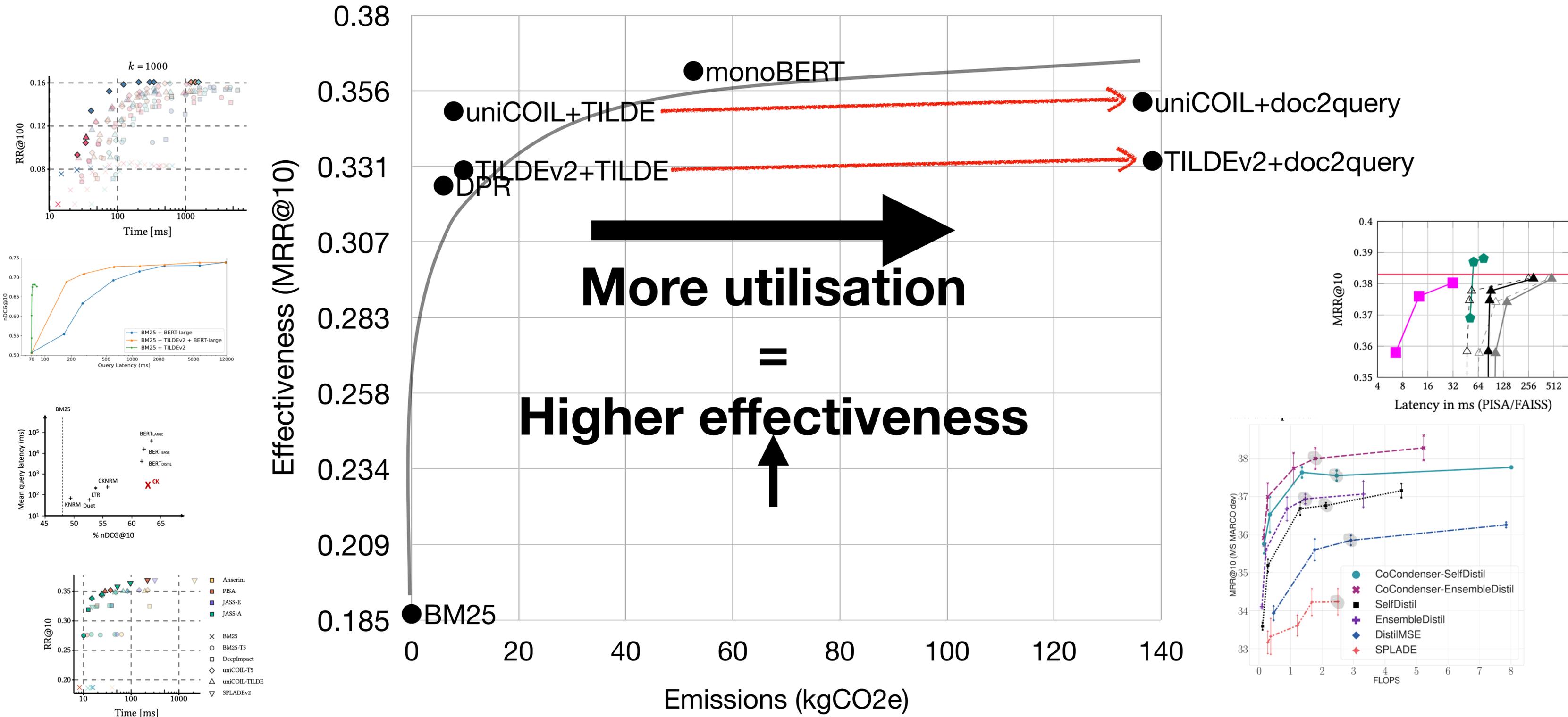
# What are the effectiveness-utilisation trade-offs of these methods?



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



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



# PART II

*Green IR in Practice*

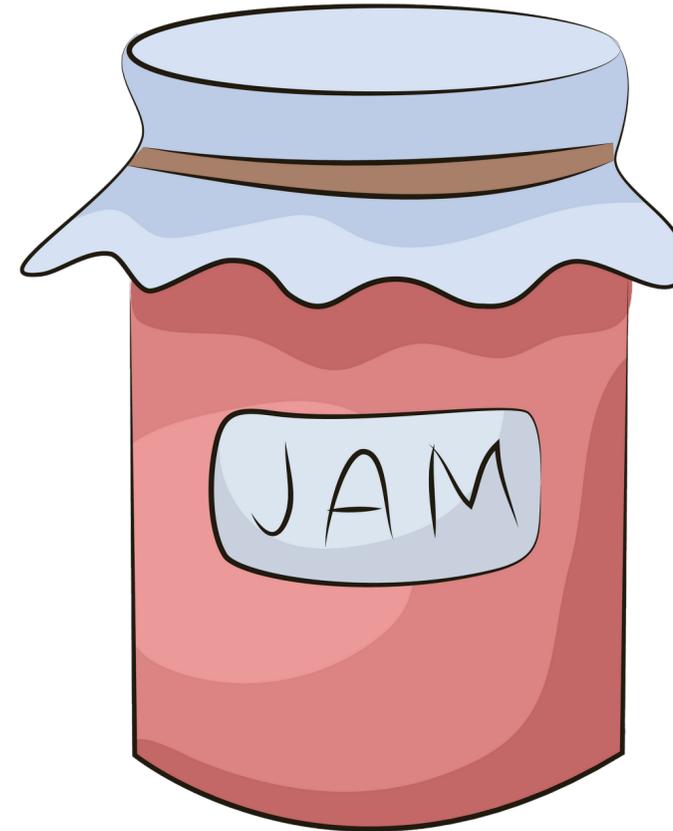


**A framework for  
practitioners to  
remain mindful of  
potential costs of  
IR research**

# Reduce



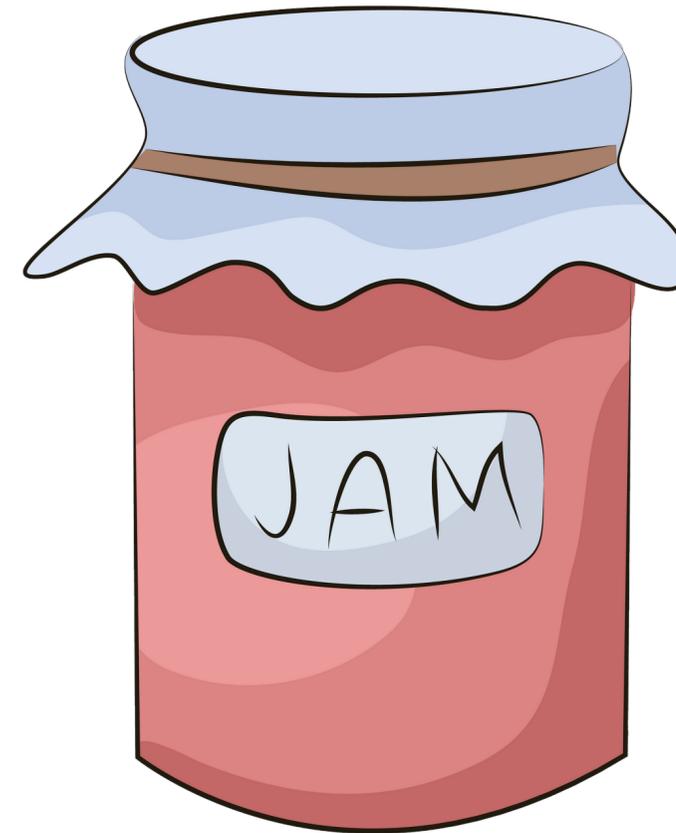
Vs



# Reduce



Vs

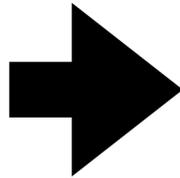


*Expend fewer resources*

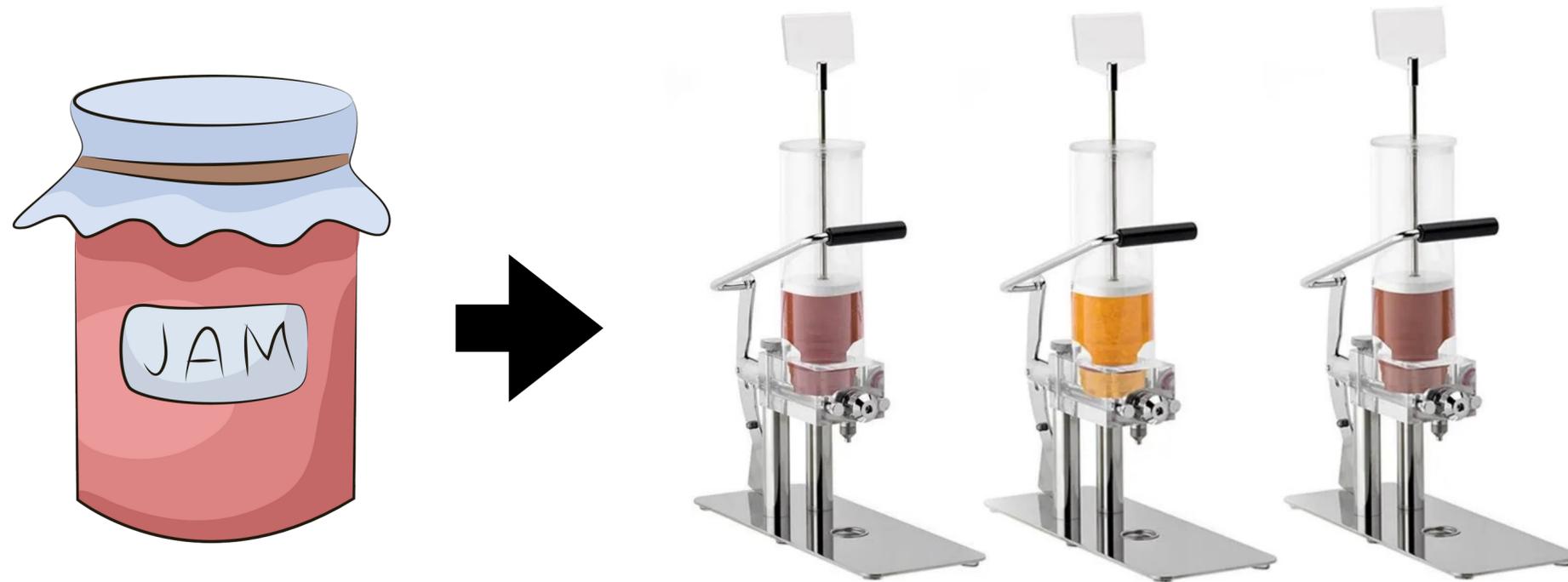
# Reduce

- 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

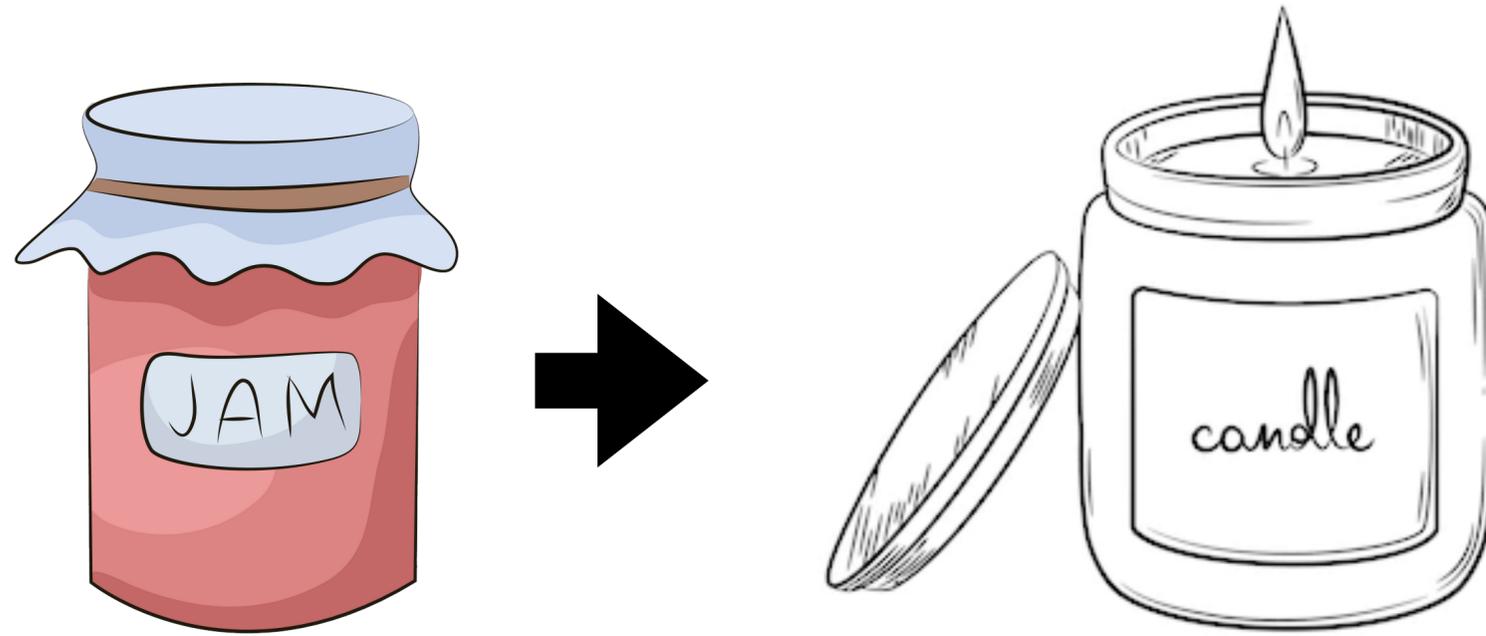


*Repurpose resources intended for one task to the same task*

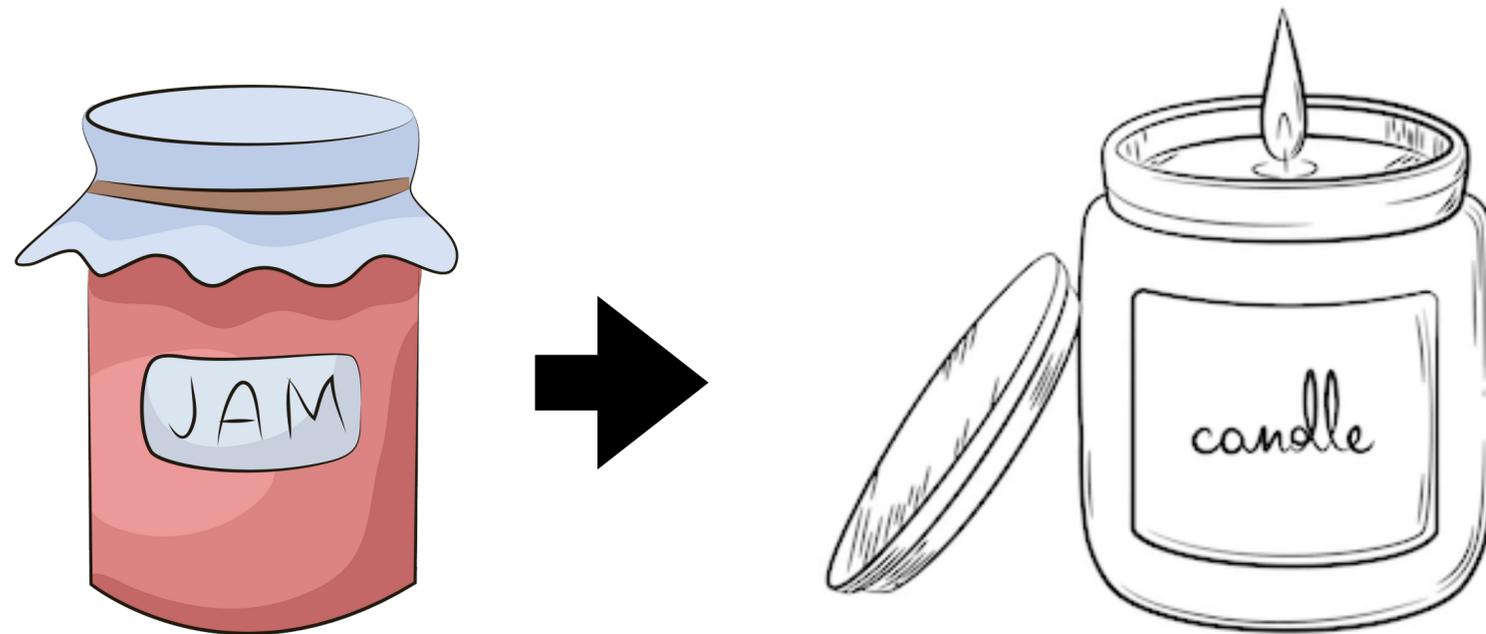
# Reuse

- Reuse existing software artefacts such as data, code, or models
- Reuse: 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 models like TILDE

# reduce, reuse, recycle

- Reduce: Expend fewer resources
- Reuse: Repurpose resources intended for one task to the same task
- Recycle: Repurpose resources intended for one task to a different task

# PART III

*Summary*

# Efficiency is not just query latency

- There is a trend of “query efficient” neural models which move the heavy computation offline
- This computation still costs: time, hardware, energy, emissions
- It is not just a “once off” cost

# Efficiency is not just latency, energy

- Data efficiency
- Learning with little data
- Frugal models, federated learning, few-shot, zero-shot, prompt learning

# Summary

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

# Summary

- **Larger neural models** = power-hungry hardware = utilisation of more power
  - However: increased 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 the entire indexing and searching architecture into a single model

# Summary

- **Larger neural models** = power-hungry hardware = utilisation of more power
  - However: increased 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 the entire indexing and searching architecture into a single model
- IR community at a **turning point**
  - Bigger/more complex models
  - Bigger collections of documents, queries

# Summary

- **Larger neural models** = power-hungry hardware = utilisation of more power
  - However: increased 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 the entire indexing and searching architecture into a single model
- IR community at a **turning point**
  - Bigger/more complex models
  - Bigger collections of documents, queries
- There is a cost to IR (+NLP, ML) research:
  - Power usage: \$\$\$
  - Emissions: CO<sub>2</sub>e