# Efficiency and Energy in Neural Information Retrieval

Harry Scells

Alexander von Humboldt Research Fellow

Leipzig University

https://scells.me

TH Köln · April 24, 2024

① Green IR

[Scells et al. 2022]

② Efficient Listwise Neural Search

[Schlatt et. al 2024]

③ Estimating Cost of IR (discussion)

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*

# Green IR
## 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 (and water [Zuccon et al. 2023])

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*

# Green IR
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 in 1 hour

- **Emissions**: by-products created by producing power
  Measured in $kgCO_2e$; 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*

# Green IR
Retrieval Efficiency

**Speed** a system can retrieve relevant information in response to a query.

Factors that can impact retrieval efficiency include:
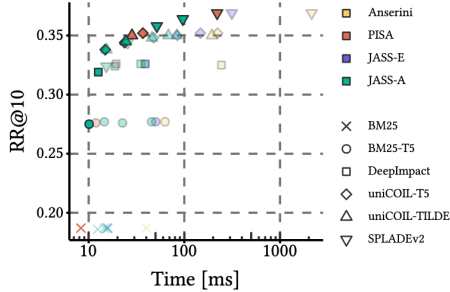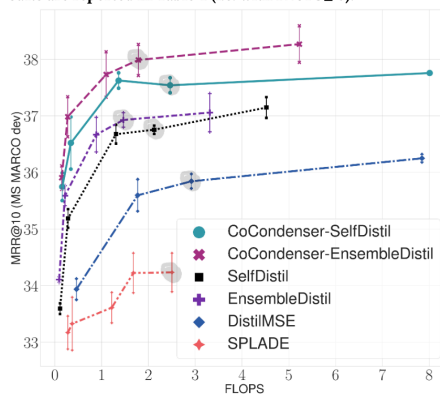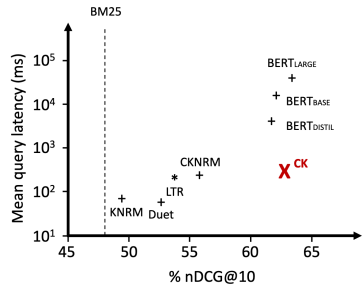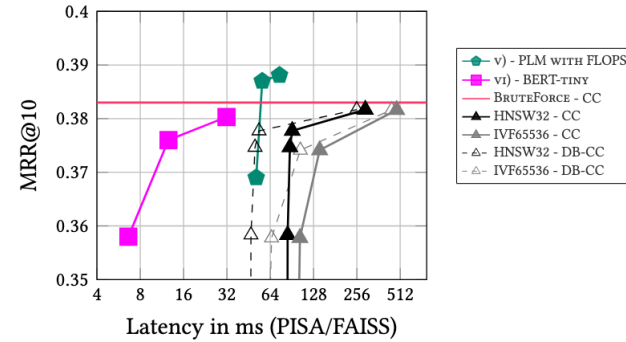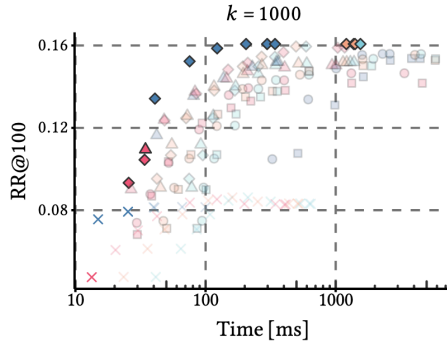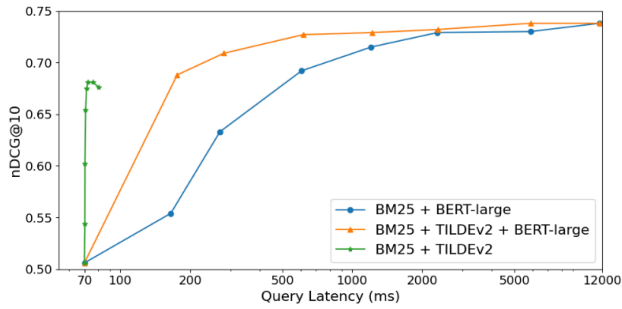
❑ **Size and complexity** of the corpus being searched

❑ Effectiveness of the **retrieval models** or techniques being used
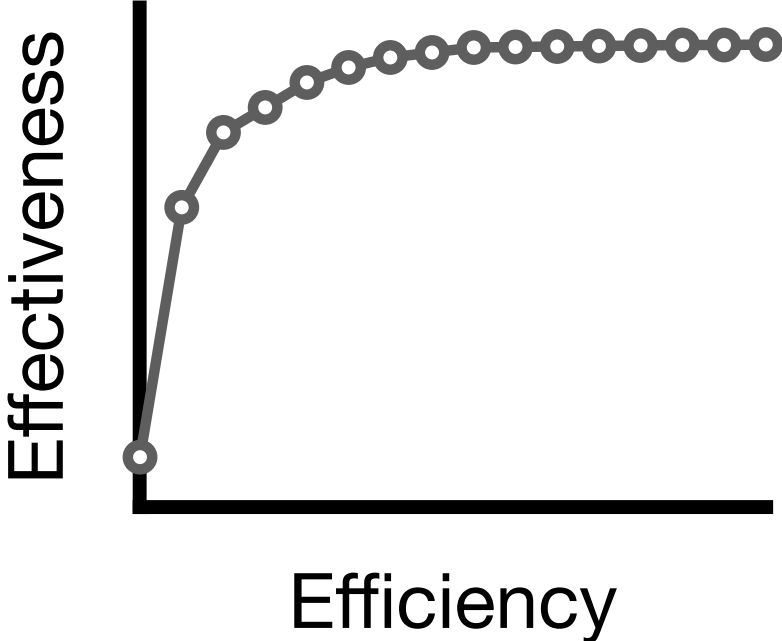
❑ Efficiency of the **hardware and infrastructure** used
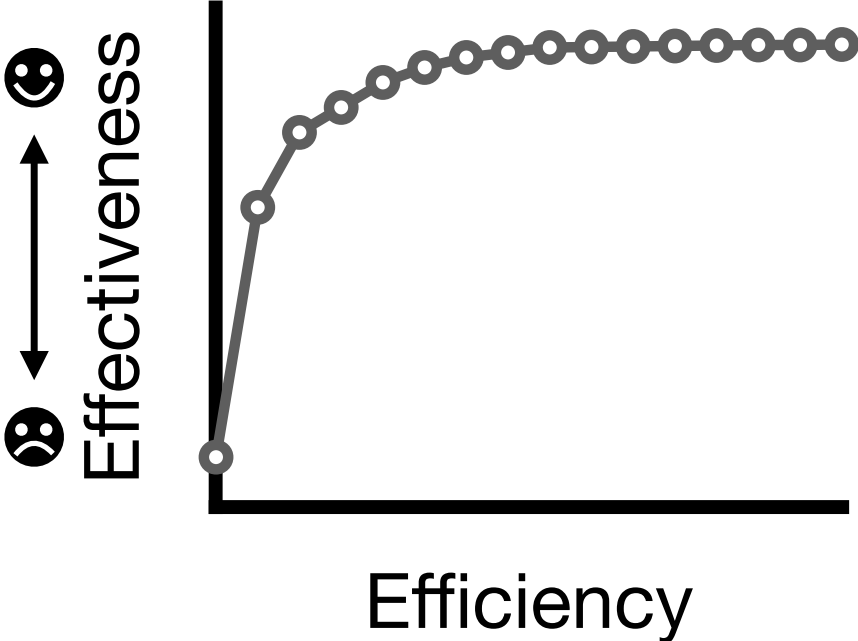
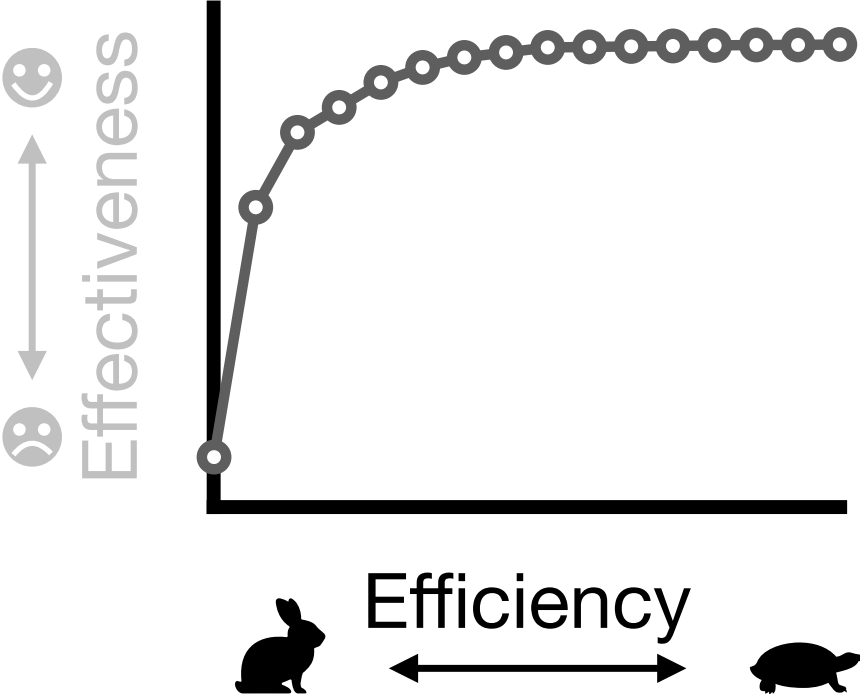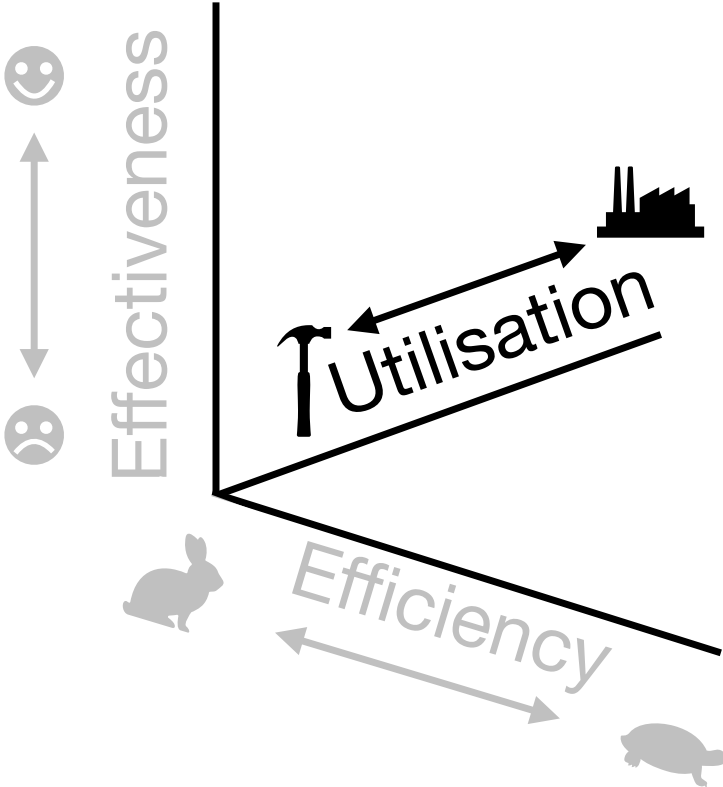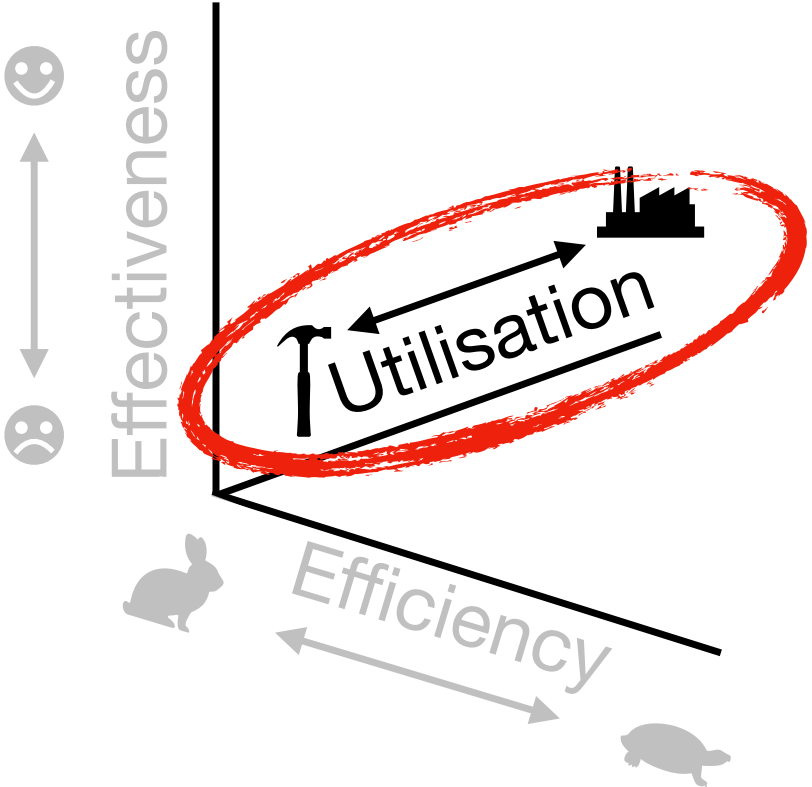# Green IR
## Retrieval Efficiency

## Retrieval Efficiency

# Green IR
## Retrieval Efficiency

## Retrieval Efficiency

Okay, so what does this mean for IR?

# Green IR

Green IR is...

*"research that yields novel results while taking into account the computational cost, encouraging a reduction in resources spent"*

(Schwartz, R. et al. 2020. Green AI. Communications of the ACM)

# Green IR

Green IR is...

*"research that yields novel results while taking into account the computational cost, encouraging a reduction in resources spent"*

(Schwartz, R. et al. 2020. Green AI. Communications of the ACM)

Neural methods require pre-trained LMs

❑ **Expensive** to create

❑ Becoming even more expensive (see: DSI and friends)

# Green 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"*

(Schwartz, R. et al. 2020. Green AI. Communications of the ACM)

Neural methods require pre-trained LMs

❑ **Expensive** to create

❑ Becoming even more expensive (see: DSI and friends)

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

# Green IR

Green IR is...

*"research that yields novel results while taking into account the computational cost, encouraging a reduction in resources spent"*

(Schwartz, R. et al. 2020. Green AI. Communications of the ACM)

Neural methods require pre-trained LMs

- ❏ **Expensive** to create

- ❏ Becoming even more expensive (see: DSI and friends)

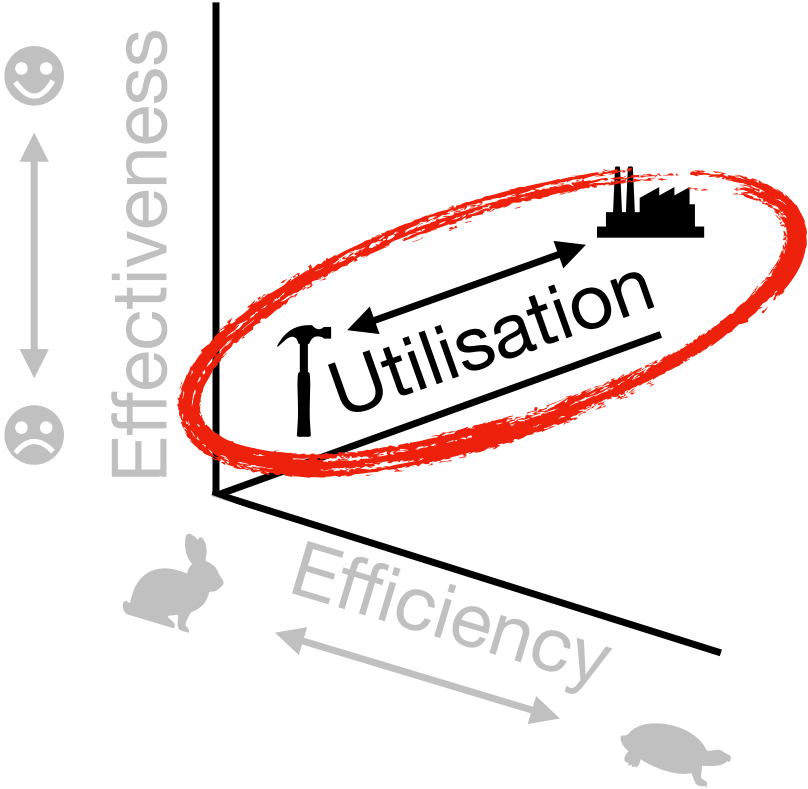**Pre-trained LMs come at a high power and emissions cost**

Missing dimension of IR evaluation

- ❏ Effectiveness
- ❏ Efficiency
- ❏ **Utilisation**

## Utilisation and Green IR

~~Okay, so what does this mean for IR?~~
Okay, so how can I measure this?

# Green IR
## Measuring Emissions

First, measure power consumption:

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

# Green IR
## Measuring Emissions

First, measure power consumption:

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

# Green IR
## Measuring Emissions

First, measure power consumption:

$$\text{PUE} \rightarrow$$

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

# Green IR
## Measuring Emissions

First, measure power consumption:

$$\underset{\textbf{watts}}{p_t} = \frac{\overset{\textbf{PUE}}{\Omega} \cdot \overset{\textbf{Running Time}}{t} \cdot (p_c + p_r + p_g)}{1000}$$

# Green IR
## Measuring Emissions

First, measure power consumption:

**PUE**  **Running Time**  **CPU, RAM, GPU power draw**

**watts**

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

# Green IR
## Measuring Emissions

First, measure power consumption:

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

**PUE** → $\Omega$

**Running Time** → $t$

**CPU, RAM, GPU power draw** → $(p_c + p_r + p_g)$

**watts** → $p_t$

Next, measure emissions:

# Green IR
## Measuring Emissions

First, measure power consumption:

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

**PUE** | **Running Time** | **CPU, RAM, GPU power draw** | **watts**

Next, measure emissions:

$$\textbf{emissions} \quad \textbf{kgCO}_2\textbf{e} = \theta \cdot p_t$$

# Green IR
## Measuring Emissions

First, measure power consumption:

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

**PUE** → $\Omega$

**Running Time** → $t$

**CPU, RAM, GPU power draw** → $(p_c + p_r + p_g)$

**watts** → $p_t$

Next, measure emissions:

$$\mathbf{kgCO_2e} = \theta \cdot p_t$$

**emissions** → $\mathbf{kgCO_2e}$

**Power consumption of experiments** → $p_t$

First, measure power consumption:

PUE    Running Time    CPU, RAM, GPU power draw

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

watts

Next, measure emissions:

avg. $CO_2$e (kg) per kWh
where experiments
took place

Power
consumption of
experiments

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

# Green IR
## Measuring Emissions

First, measure power consumption:

PUE      **Running Time**      **CPU, RAM, GPU power draw**

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

**watts**

Next, measure emissions:

avg. $CO_2$e (kg) per kWh
where experiments
took place

**Power
consumption of
experiments**

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

Emissions of my search engine:

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

# Green IR
## Measuring Emissions

First, measure power consumption:

PUE    **Running Time**    **CPU, RAM, GPU power draw**

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

**watts** $\rightarrow p_t$

Next, measure emissions:

**avg. $CO_2e$ (kg) per kWh where experiments took place**

**Power consumption of experiments**

emissions $\rightarrow$ $$\mathbf{kgCO_2e} = \theta \cdot p_t$$

Emissions of my search engine:

$$\mathbf{kgCO_2e} = \theta \cdot \Delta_q \cdot p_q$$

**Power consumption of a single query**

# Green IR
## Measuring Emissions

First, measure power consumption:

PUE     **Running Time**     **CPU, RAM, GPU power draw**

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

**watts**

Next, measure emissions:

**avg. $CO_2e$ (kg) per kWh where experiments took place**

**Power consumption of experiments**

$$\text{emissions} \rightarrow \mathbf{kgCO_2e} = \theta \cdot p_t$$

Emissions of my search engine:

**No. queries issued per unit time**

**Power consumption of a single query**

$$\mathbf{kgCO_2e} = \theta \cdot \Delta_q \cdot p_q$$

## Utilisation and Green IR



~~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!

# Green IR

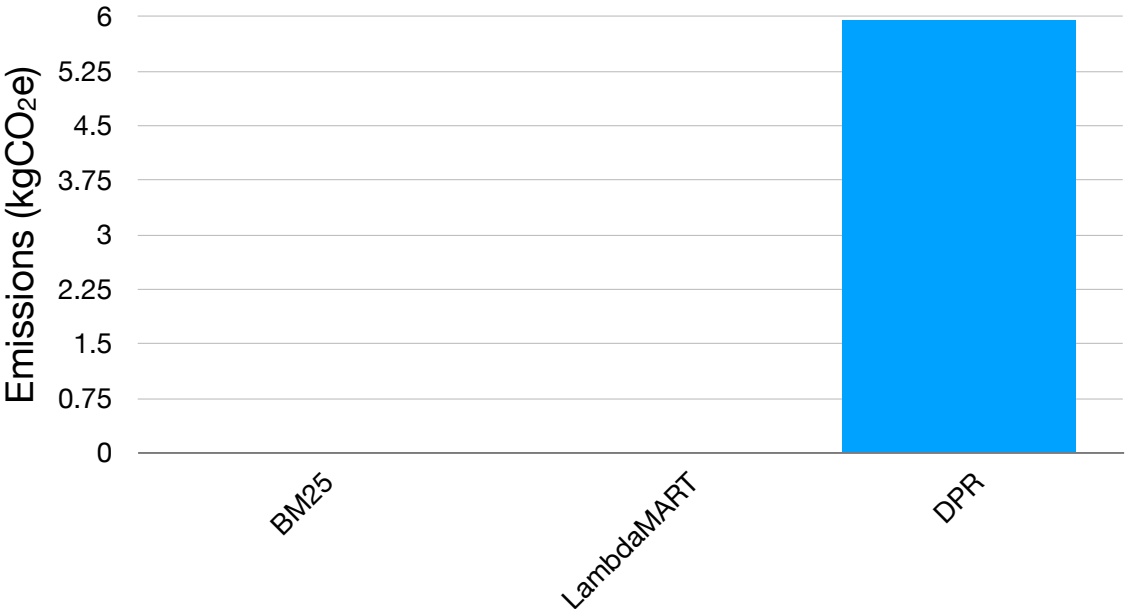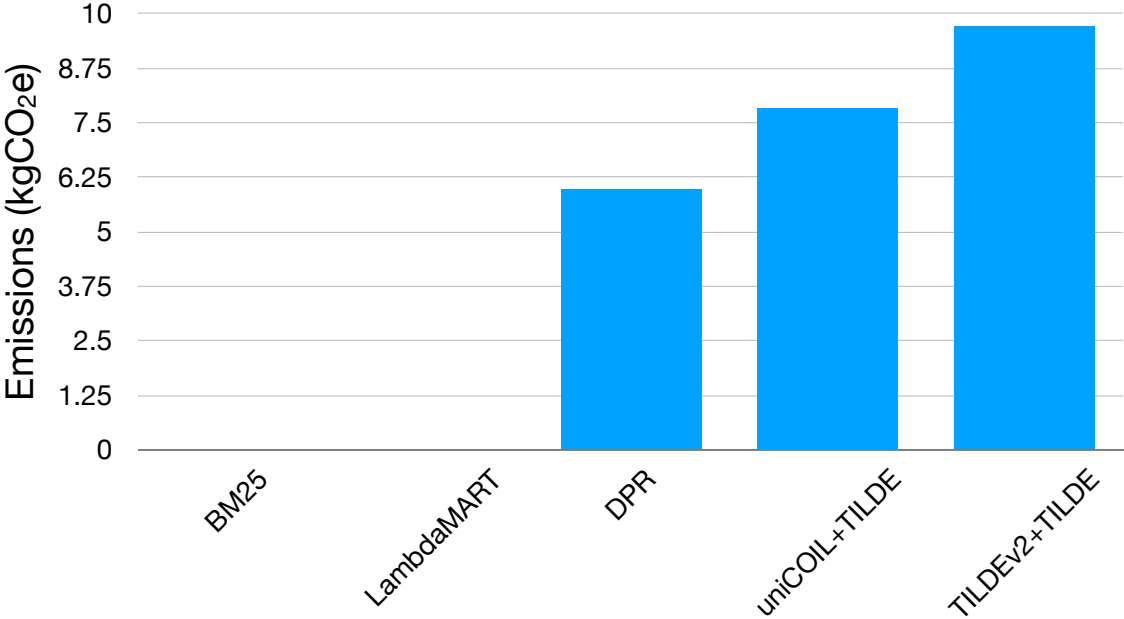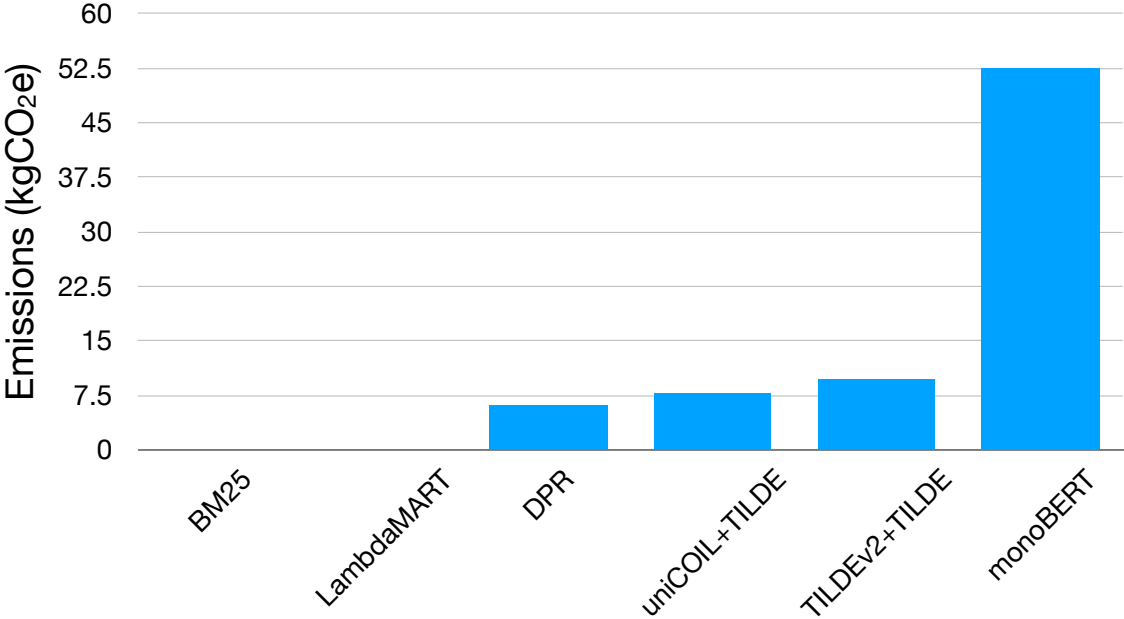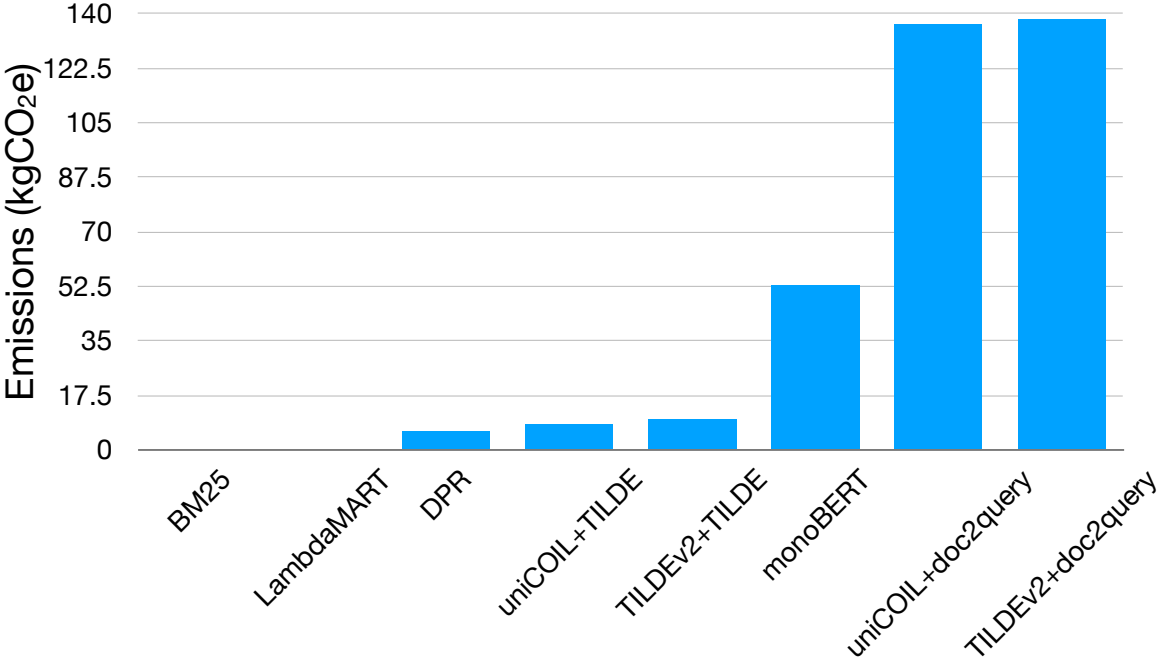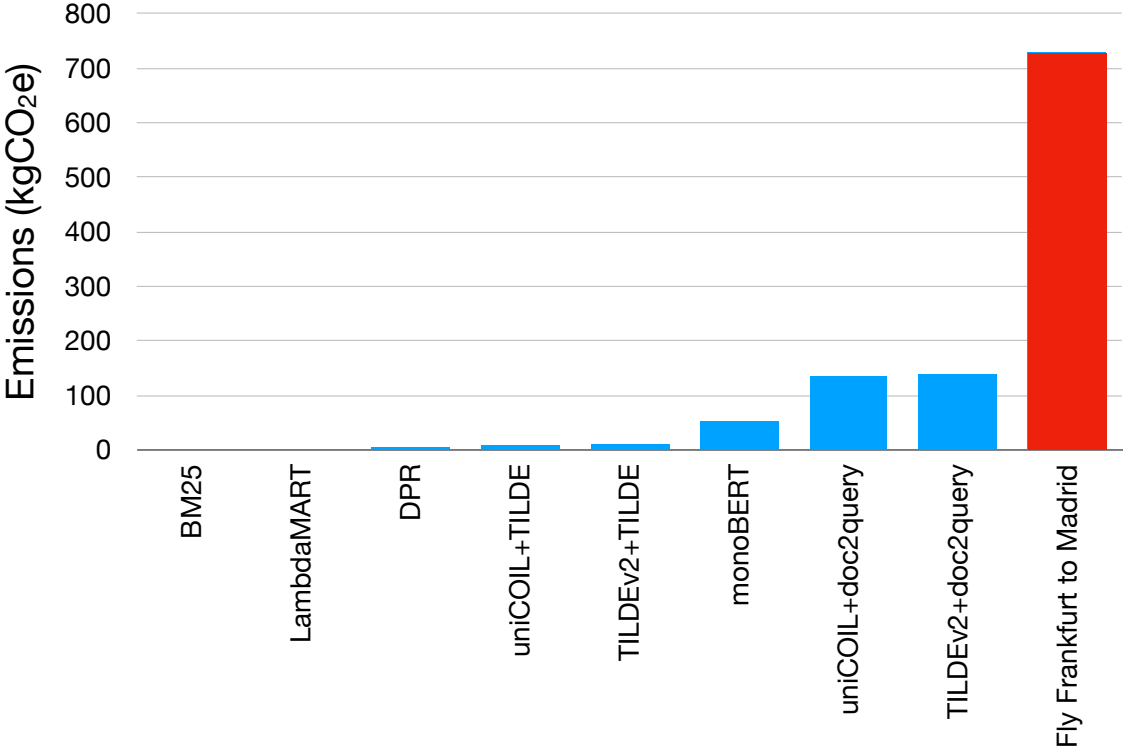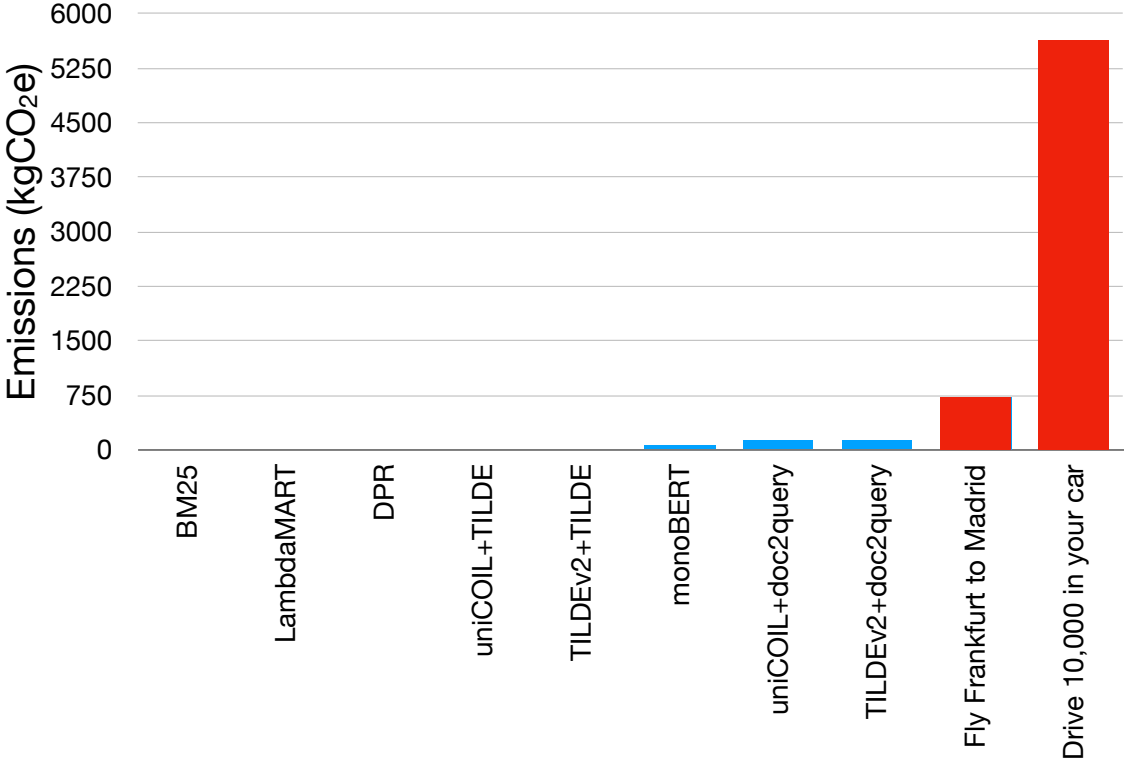## How many emissions produced to obtain a single result?

# Green IR

How many emissions produced to obtain a single result?

# Green IR

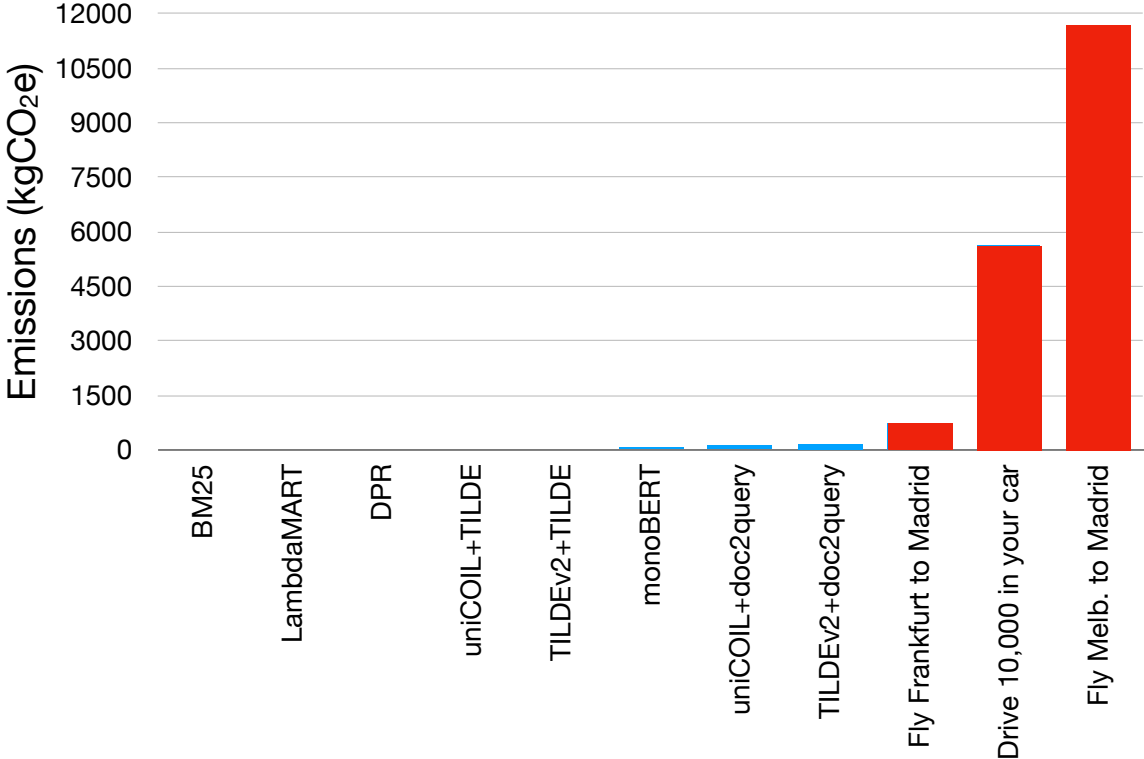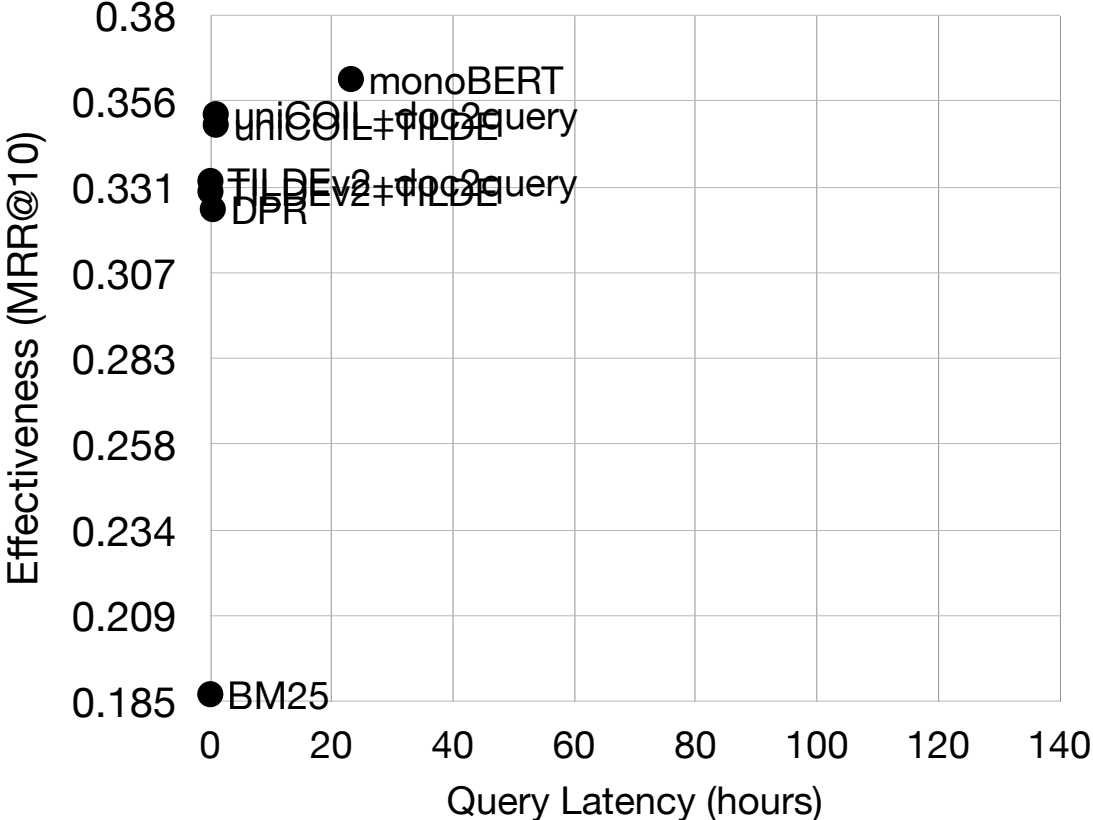How many emissions produced to obtain a single result?

# Green IR

How many emissions produced to obtain a single result?
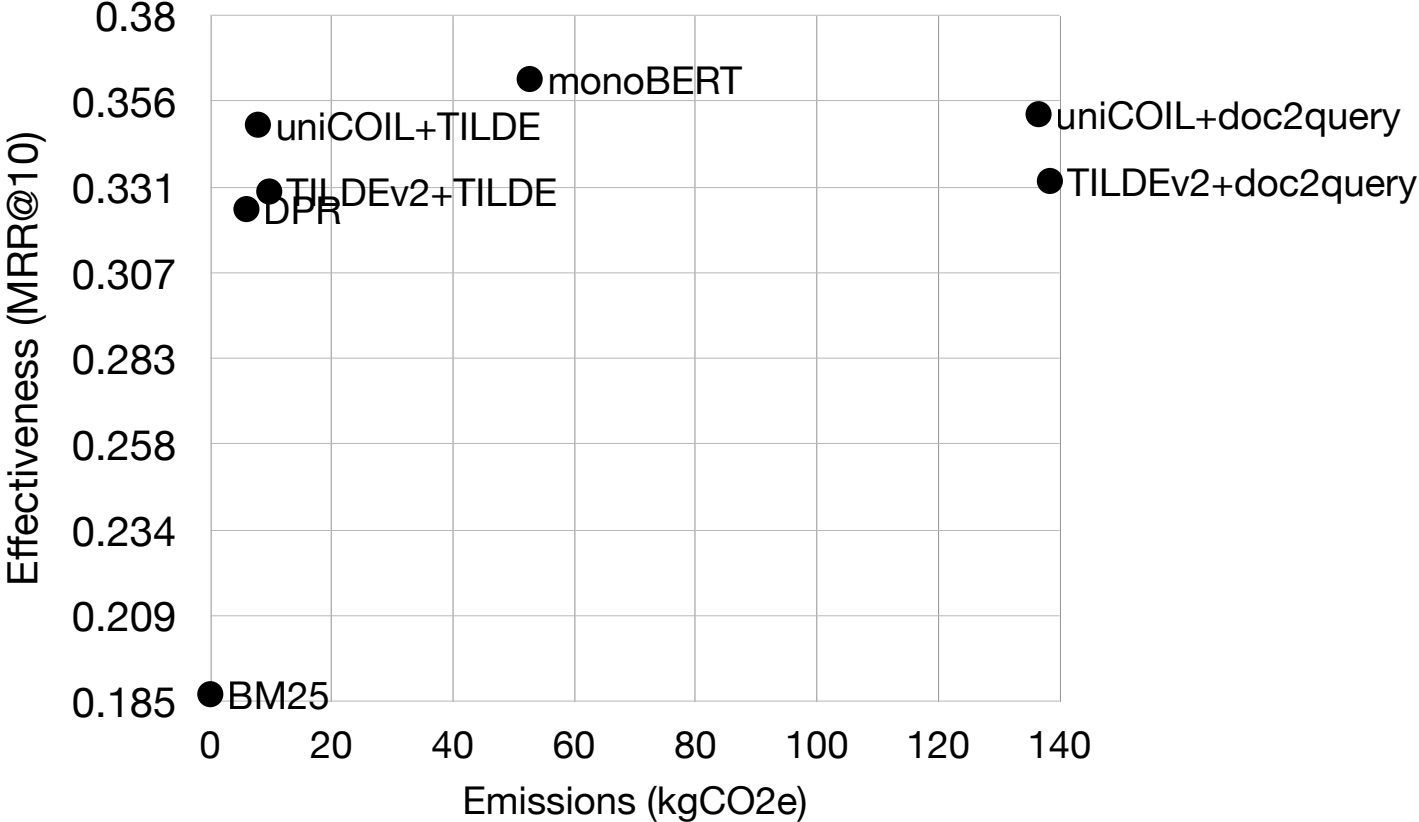


Bar chart titled "Emissions (kgCO₂e)" on the y-axis (0 to 60, in increments of 7.5) versus model names on the x-axis. Values approximately: BM25 ≈ 0, LambdaMART ≈ 0, DPR ≈ 6, uniCOIL+TILDE ≈ 8, TILDEv2+TILDE ≈ 9.5, monoBERT ≈ 52.5.

# Green IR

How many emissions produced to obtain a single result?

# Green IR

How many emissions produced to obtain a single result?



Bar chart — Emissions ($kgCO_2e$) by method:
- BM25: ~0
- LambdaMART: ~0
- DPR: ~5
- uniCOIL+TILDE: ~8
- TILDEv2+TILDE: ~10
- monoBERT: ~50
- uniCOIL+doc2query: ~135
- TILDEv2+doc2query: ~140
- Fly Frankfurt to Madrid: ~730 (red)

# Green IR

How many emissions produced to obtain a single result?



Bar chart titled with y-axis "Emissions (kgCO₂e)" ranging from 0 to 6000. Categories: BM25, LambdaMART, DPR, uniCOIL+TILDE, TILDEv2+TILDE, monoBERT, uniCOIL+doc2query, TILDEv2+doc2query, Fly Frankfurt to Madrid (~750), Drive 10,000 in your car (~5600).

# Green IR

How many emissions produced to obtain a single result?



Bar chart titled with y-axis "Emissions (kgCO₂e)" ranging from 0 to 12000. Categories along the x-axis: BM25, LambdaMART, DPR, uniCOIL+TILDE, TILDEv2+TILDE, monoBERT, uniCOIL+doc2query, TILDEv2+doc2query, Fly Frankfurt to Madrid, Drive 10,000 in your car, Fly Melb. to Madrid.
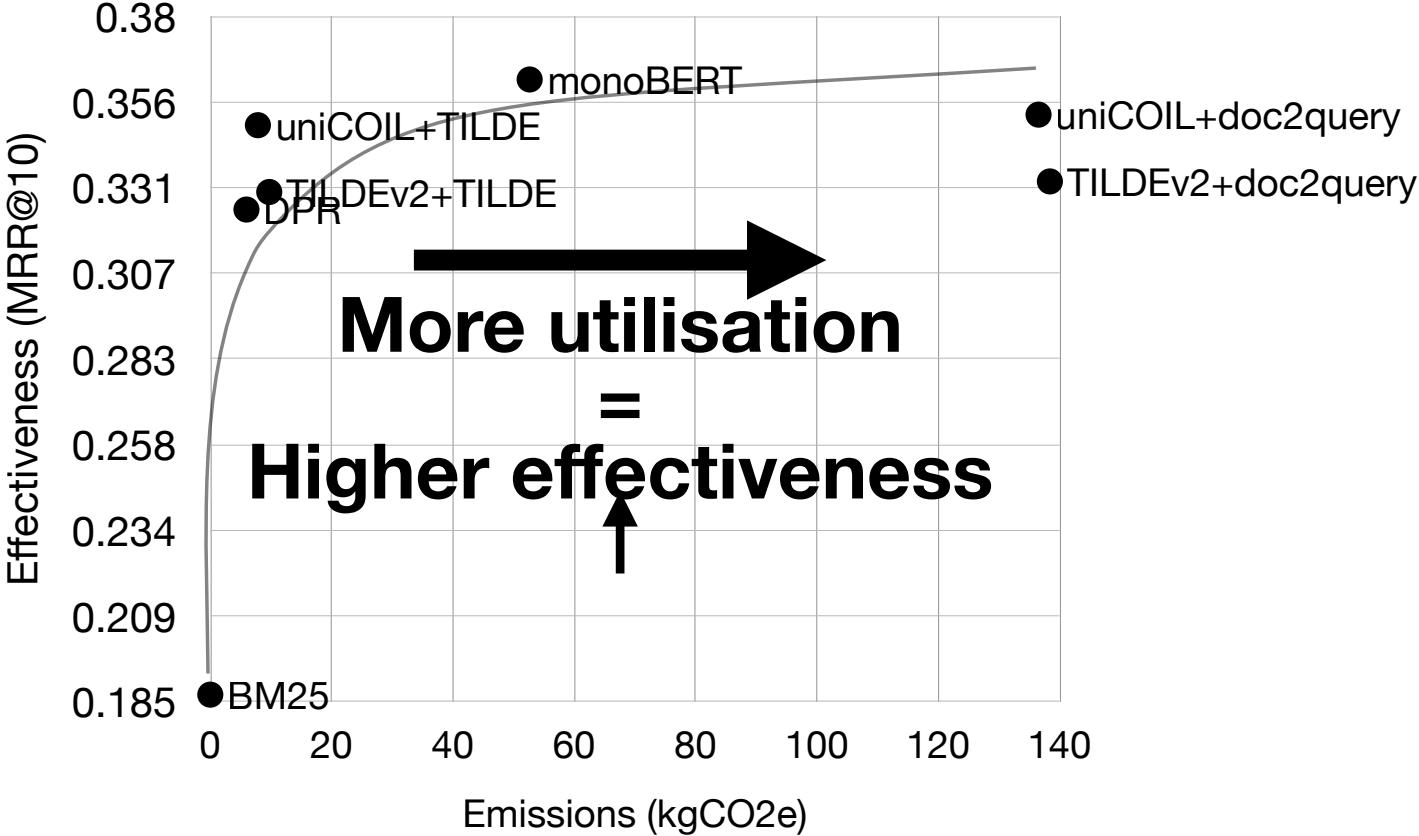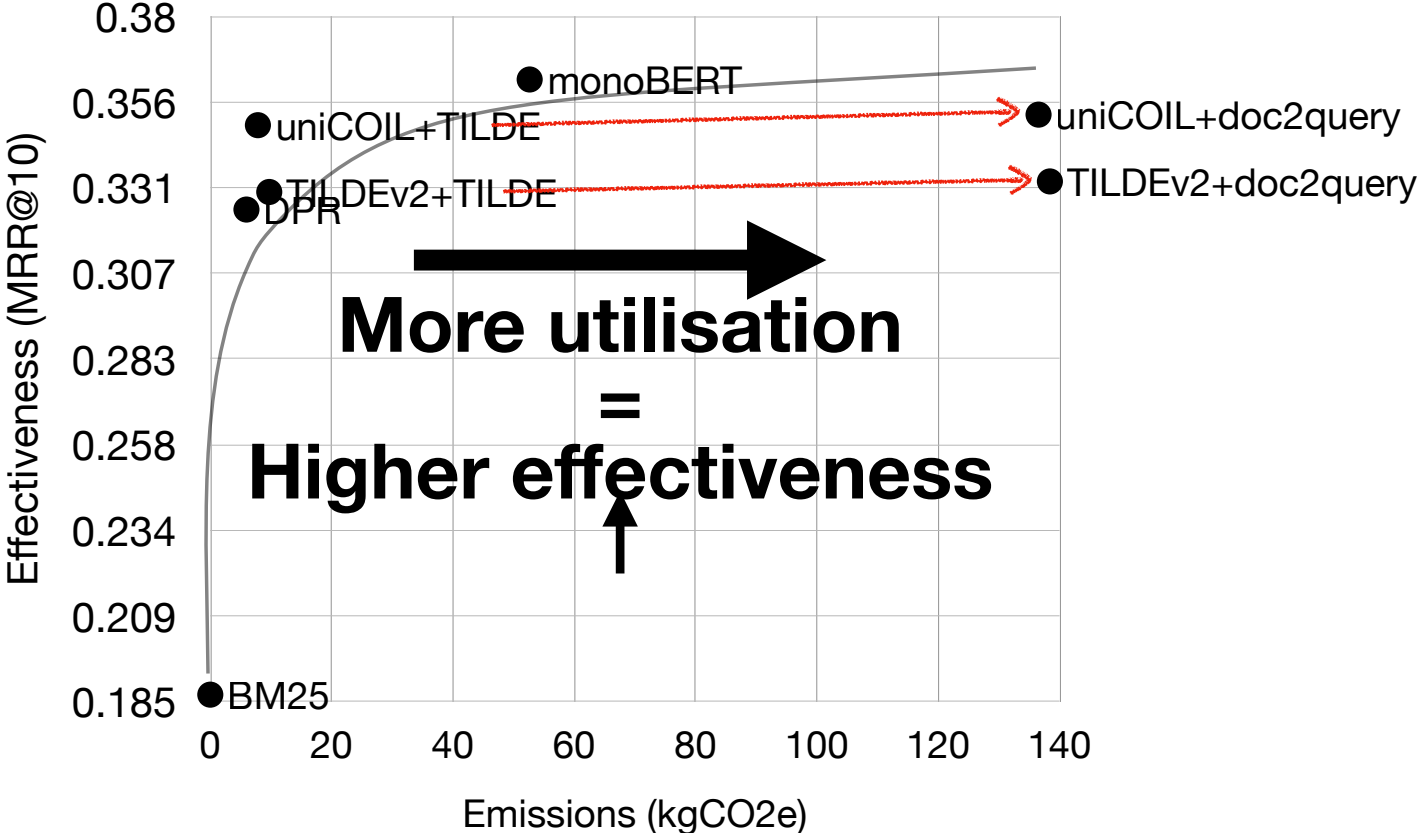
# Green IR

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

# Green IR

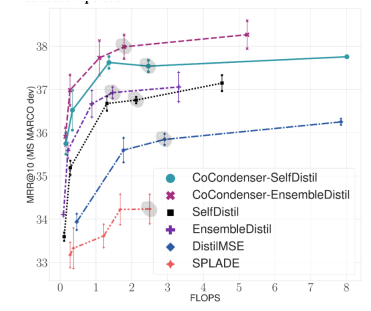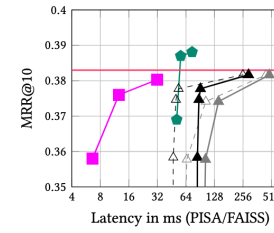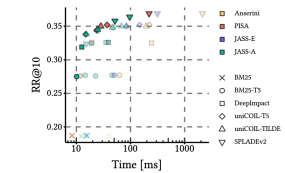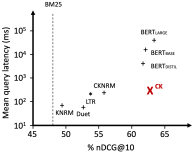What are the effectiveness-utilisation trade-offs of these methods?



Scatter plot. Y-axis: Effectiveness (MRR@10), ranging from 0.185 to 0.38. X-axis: Emissions (kgCO2e), ranging from 0 to 140.

- BM25: ~(0, 0.185)
- DPR: ~(6, 0.325)
- TILDEv2+TILDE: ~(9, 0.331)
- uniCOIL+TILDE: ~(8, 0.350)
- monoBERT: ~(52, 0.362)
- TILDEv2+doc2query: ~(137, 0.331)
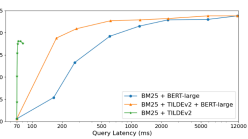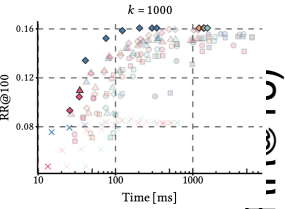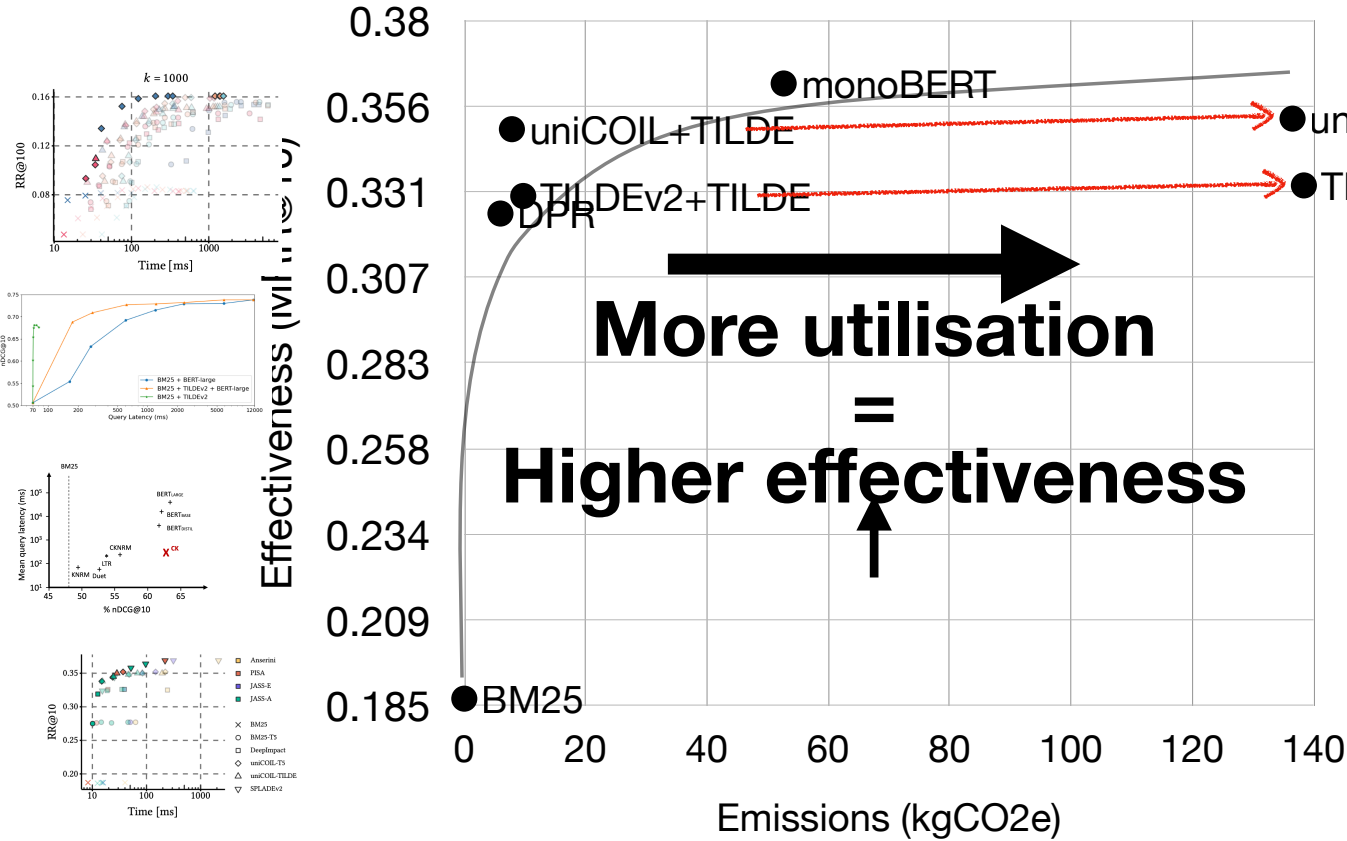- uniCOIL+doc2query: ~(137, 0.352)

# Green IR

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

# Green IR

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

# Green IR

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

# Green IR

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

# Green IR

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



**More utilisation = Higher effectiveness**

Plot axes: Effectiveness (MRR@10?) vs Emissions (kgCO2e)

Data points: BM25 (0.185), DPR (0.331), TILDEv2+TILDE (0.331), uniCOIL+TILDE (0.356), monoBERT (~0.365), uniCOIL+doc2query (0.356), TILDEv2+doc2query (0.331)

# Green IR
Reduce, Reuse, Recycle

Reduce ➜ expend fewer resources

- ❏ 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?

# Green IR
## Reduce, Reuse, Recycle

Reduce ➜ expend fewer resources

- ❏ 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?

Reuse ➜ repurpose resources intended for one task to the same task

- ❏ 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 or code meant for one task to the same task?

# Green IR

Reduce, Reuse, Recycle

Reduce ➡ expend fewer resources

❑ 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?

Reuse ➡ repurpose resources intended for one task to the same task

❑ 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 or code meant for one task to the same task?

Recycle ➡ repurpose resources intended for one task to a different task

❑ Recycle existing software artefacts such as data, code, or models

❑ 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 or code meant for one task to a different task?

① Green IR

[Scells et al. 2022]

② Efficient Listwise Neural Search

[Schlatt et. al 2024]

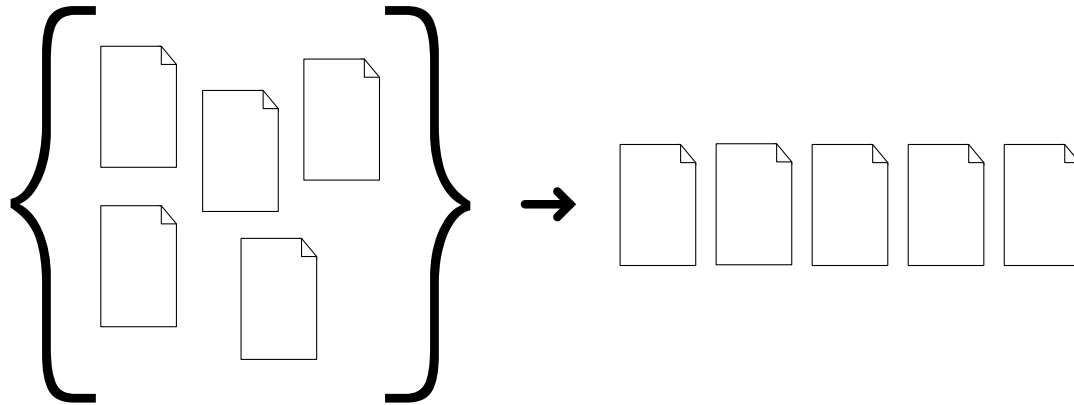③ Estimating Cost of IR (discussion)

# Efficient Listwise Neural Search

## Motivation [Schlatt et. al 2024]

**Learning task:** Given a set of objects, rank them according to a ranking criterion

- ❑ Ranking of documents from a set of documents and a query

- ❑ Existing transformer architecture cannot model this task effectively

- ❑ Two properties: Permutation invariance and cross-document information

# Efficient Listwise Neural Search
## Motivation [Schlatt et. al 2024]

**Learning task:** Given a set of objects, rank them according to a ranking criterion

- ❏ Ranking of documents from a set of documents and a query

- ❏ Existing transformer architecture cannot model this task effectively

- ❏ Two properties: Permutation invariance and cross-document information



Existing architectures model either one of these properties **but never both**

- ❏ Trade off effective ranking for permutation invariance ➜ Pointwise

- ❏ Trade off efficient ranking for cross-document information ➜ Listwise

# Efficient Listwise Neural Search
## Model Architecture

Pointwise

- More efficient at the expense of effectiveness.

- Permutation-invariant, no cross-document information.

- Scaleable: each query-document pair is scored.

**Listwise**

- More effective at the expense of efficiency.

- Non-permutation-invariant, cross-document information.

- Unscaleable: All permutations of query-documents is scored.

# Efficient Listwise Neural Search
## Model Architecture

Pointwise

- More efficient at the expense of effectiveness.

- Permutation-invariant, no cross-document information.

- Scaleable: each query-document pair is scored.

## **Listwise**

- More effective at the expense of efficiency.

- Non-permutation-invariant, cross-document information.

- Unscaleable: All permutations of query-documents is scored.

State of the Art





Batch

CLS token    Query tokens    Document tokens

# Efficient Listwise Neural Search
## Model Architecture

Document scoring:

- ❑ Each permutation of documents is fed into model.

- ❑ Reason: Transformer is sequence modeller; order of documents biases the score.



State of the Art

# Efficient Listwise Neural Search
## Model Architecture

Document scoring:

- ❏ Each permutation of documents is fed into model.

- ❏ Reason: Transformer is sequence modeller; order of documents biases the score.

- ❏ Task: Predict ordering preference of documents given query.



State of the Art

Query

Doc 1

Doc 2

Doc 3

P1 > P2 > P3
P2 > P1 > P3
P1 > P3 > P2
P1 > P2 > P3
P1 > P2 > P3
P1 > P3 > P2

Batch

CLS token    Query tokens    Document tokens

# Efficient Listwise Neural Search

## Model Architecture

Document scoring:

- ❏ Each permutation of documents is fed into model.

- ❏ Reason: Transformer is sequence modeller; order of documents biases the score.

- ❏ Task: Predict ordering preference of documents given query.

- ❏ Score computed by aggregating preferences.



State of the Art

Query

Doc 1    0.9
Doc 2    0.5
Doc 3    0.1

P1 > P2 > P3
P2 > P1 > P3
P1 > P3 > P2
P1 > P2 > P3
P1 > P2 > P3
P1 > P3 > P2

Batch

CLS token    Query tokens    Document tokens

# Efficient Listwise Neural Search
## Model Architecture

Set-Encoder document scoring:

❑ Each query-document pair only needs to be scored once.

# Efficient Listwise Neural Search
## Model Architecture

Set-Encoder document scoring:

❑ Each query-document pair only needs to be scored once.

❑ Share cross-document information through attention mechanism.

❑ Reset positional information to make scores permutation invariant.

# Efficient Listwise Neural Search
## Model Architecture

Set-Encoder document scoring:

- ❏ Each query-document pair only needs to be scored once.

- ❏ Share cross-document information through attention mechanism.

- ❏ Reset positional information to make scores permutation invariant.

- ❏ Score computed directly for all query-document pairs.

# Efficient Listwise Neural Search

## Modelling Cross-Document Interactions with Attention

q, d$_1$

q, d$_2$

$$\text{Attention}(\mathsf{Q}, \mathsf{K}, \mathsf{V}) = \text{softmax}(\tfrac{\mathsf{Q}\mathsf{K}^T}{\sqrt{h}})V$$

# Efficient Listwise Neural Search

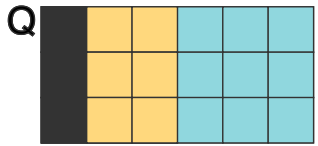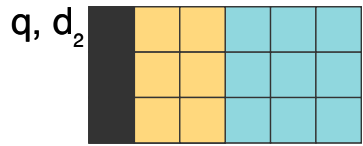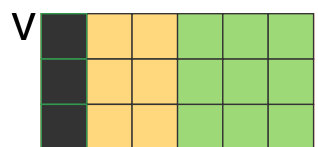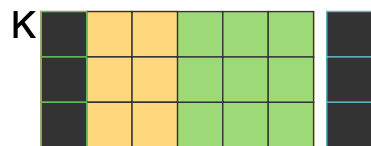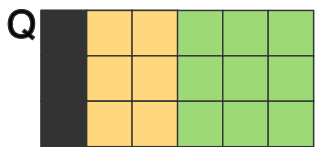## Modelling Cross-Document Interactions with Attention

q, $d_1$

Q

K

V

q, $d_2$

Q

K

V
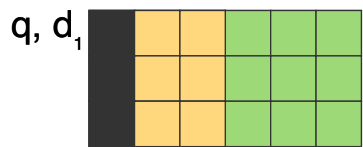
$$\text{Attention}(Q, K, V) = \text{softmax}(\frac{QK^T}{\sqrt{h}})V$$

# Efficient Listwise Neural Search
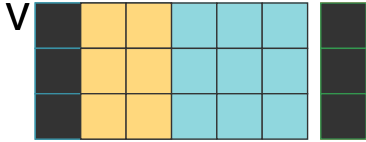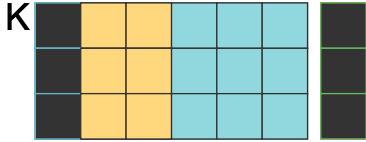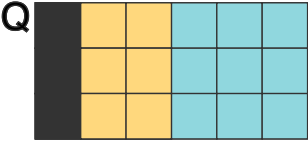
## Modelling Cross-Document Interactions with Attention

q, d₁

Q

K

V

q, d₂

Q

K

V

$$\text{Attention}(Q, K, V) = \text{softmax}(\frac{QK^T}{\sqrt{h}})V$$

For $d_i$, let $\bar{K}^i = [K_1^j : j \neq i]$

# Efficient Listwise Neural Search
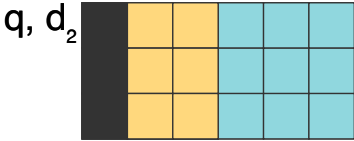
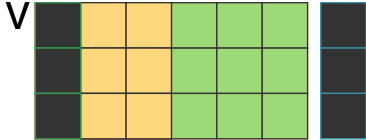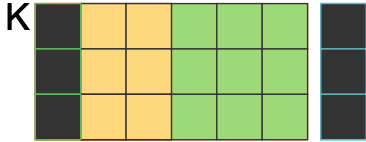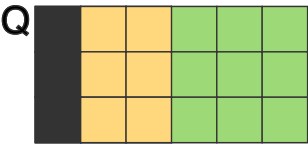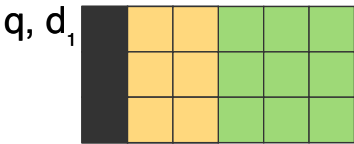## Modelling Cross-Document Interactions with Attention



$$\text{Attention}(Q, K, V) = \text{softmax}\left(\frac{QK^T}{\sqrt{h}}\right)V$$

For $d_i$, let $\bar{K}^i = [K_1^j : j \neq i]$

For $d_i$, let $\bar{V}^i = [V_1^j : j \neq i]$

# Efficient Listwise Neural Search

## Modelling Cross-Document Interactions with Attention



$$\text{Attention}(\mathsf{Q}, \mathsf{K}, \mathsf{V}) = \text{softmax}\left(\frac{\mathsf{Q}\mathsf{K}^T}{\sqrt{h}}\right)V$$

For $d_i$, let $\bar{K}^i = [K_1^j : j \neq i]$
For $d_i$, let $\bar{V}^i = [V_1^j : j \neq i]$

Cross-document attention for $d_i$:
$\text{Attention}(\mathsf{Q}^i, [\mathsf{K}^i\bar{\mathsf{K}}^i], [\mathsf{V}^i\bar{\mathsf{V}}^i])$

# Efficient Listwise Neural Search

## Attention Visualised

# Efficient Listwise Neural Search

## Attention Visualised



**Typical Cross-Encoder**

| Token Type | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| [CLS] | 0.28 | 0.15 | 0.26 | 0.21 | 0.22 | 0.27 | 0.26 | 0.18 | 0.06 | 0.02 | 0.05 | 0.02 |
| Query | 0.06 | 0.15 | 0.27 | 0.23 | 0.25 | 0.28 | 0.27 | 0.22 | 0.11 | 0.08 | 0.15 | 0.53 |
| Document | 0.66 | 0.70 | 0.47 | 0.55 | 0.52 | 0.44 | 0.47 | 0.59 | 0.83 | 0.90 | 0.80 | 0.45 |
| Other [CLS] | | | | | | | | | | | | |

**Set-Encoder**

| Token Type | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| [CLS] | 0.01 | 0.01 | 0.01 | 0.01 | 0.01 | 0.01 | 0.01 | 0.01 | 0.01 | 0.01 | 0.01 | 0.00 |
| Query | 0.02 | 0.01 | 0.01 | 0.01 | 0.01 | 0.00 | 0.00 | 0.01 | 0.03 | 0.02 | 0.09 | 0.38 |
| Document | 0.20 | 0.16 | 0.04 | 0.07 | 0.08 | 0.02 | 0.02 | 0.04 | 0.22 | 0.39 | 0.29 | 0.37 |
| Other [CLS] | 0.77 | 0.82 | 0.94 | 0.92 | 0.90 | 0.97 | 0.97 | 0.95 | 0.74 | 0.58 | 0.61 | 0.25 |

Layer

**Set-Encoder attends to other documents in early layers, then the document to score in final layers.**

# Efficient Listwise Neural Search

## Results: Ranking Effectiveness

| Model | Parameters | Effectiveness (nDCG@10) |
|---|---|---|
| monoBERT base | 110M | 0.379 |
| monoBERT large | 340M | 0.381 |
| monoT5 base | 220M | 0.376 |
| monoT5 large | 3B | 0.410 |
| LiT5-Distill | 220M | 0.406 |
| Set-Encoder | 110M | 0.406 |

# Efficient Listwise Neural Search

## Results: Ranking Effectiveness

| Model | Parameters | Effectiveness (nDCG@10) |
|---|---|---|
| monoBERT base | 110M | 0.379 |
| monoBERT large | 340M | 0.381 |
| monoT5 base | 220M | 0.376 |
| monoT5 large | 3B | 0.410 |
| LiT5-Distill | 220M | 0.406 |
| Set-Encoder | 110M | 0.406 |

**Set-Encoder has same effectiveness of SOTA listwise model with half the parameters.**

# Efficient Listwise Neural Search

## Results: Ranking Effectiveness

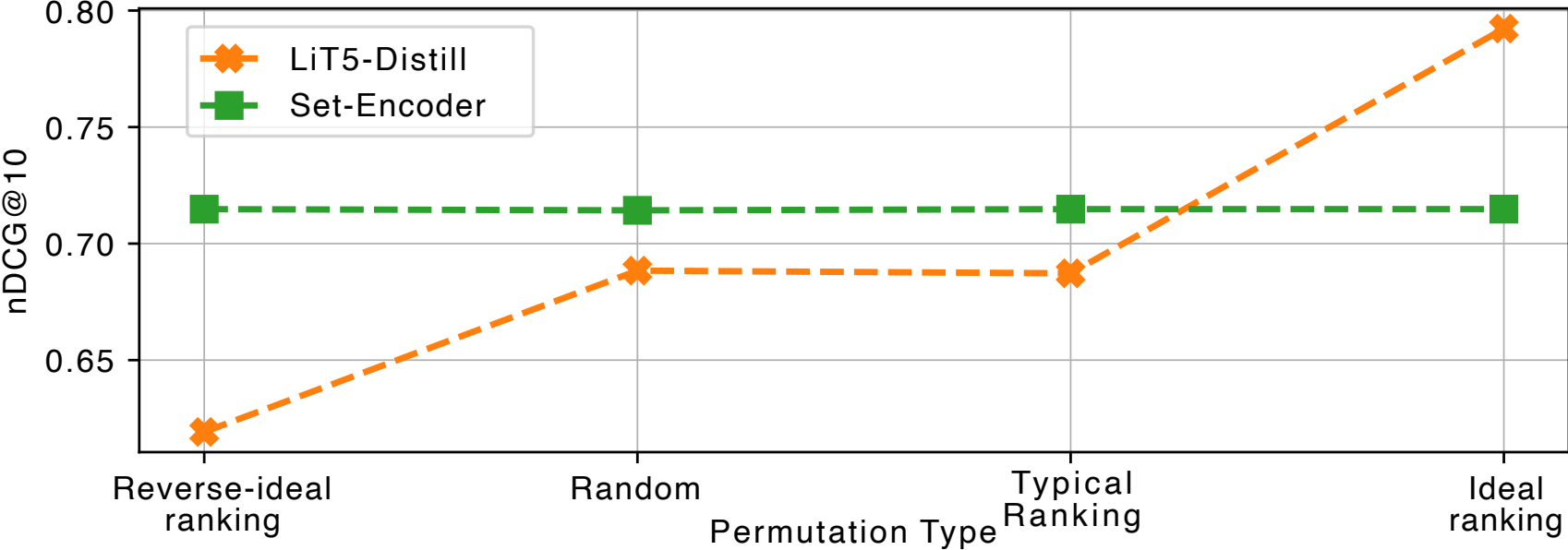| Model | Parameters | Effectiveness (nDCG@10) |
|---|---|---|
| monoBERT base | 110M | 0.379 |
| monoBERT large | 340M | 0.381 |
| monoT5 base | 220M | 0.376 |
| monoT5 large | 3B | 0.410 |
| LiT5-Distill | 220M | 0.406 |
| Set-Encoder | 110M | 0.406 |

**Set-Encoder has same effectiveness of SOTA listwise model with half the parameters.**

**Set-Encoder has similar effectiveness to SOTA pointwise model with 3B fewer parameters.**
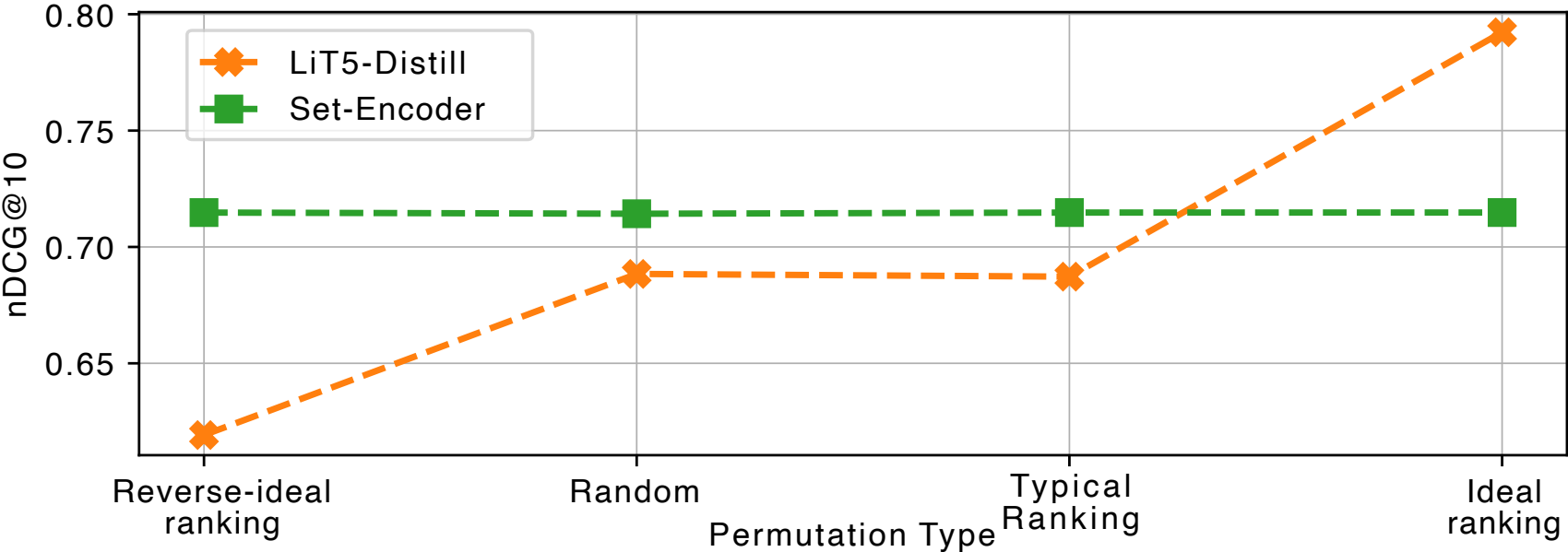
# Efficient Listwise Neural Search

## Robustness to Initial Ranking Permutations

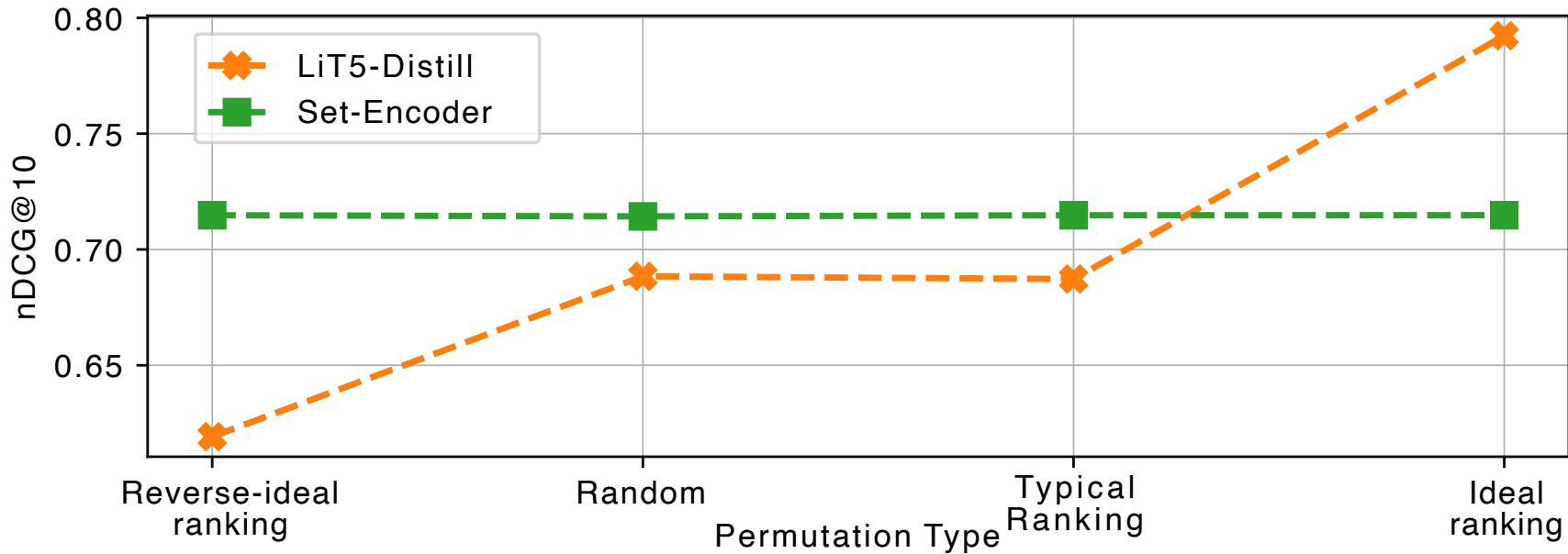# Efficient Listwise Neural Search
## Robustness to Initial Ranking Permutations



Irrespective of initial document ranking,
Set-Encoder has same effectiveness.

# Efficient Listwise Neural Search
## Robustness to Initial Ranking Permutations



**Irrespective of initial document ranking, Set-Encoder has same effectiveness.**

**SOTA Listwise model makes document ranking worse when given ideal ranking.**

(1) Green IR

[Scells et al. 2022]

(2) Efficient Listwise Neural Search

[Schlatt et. al 2024]

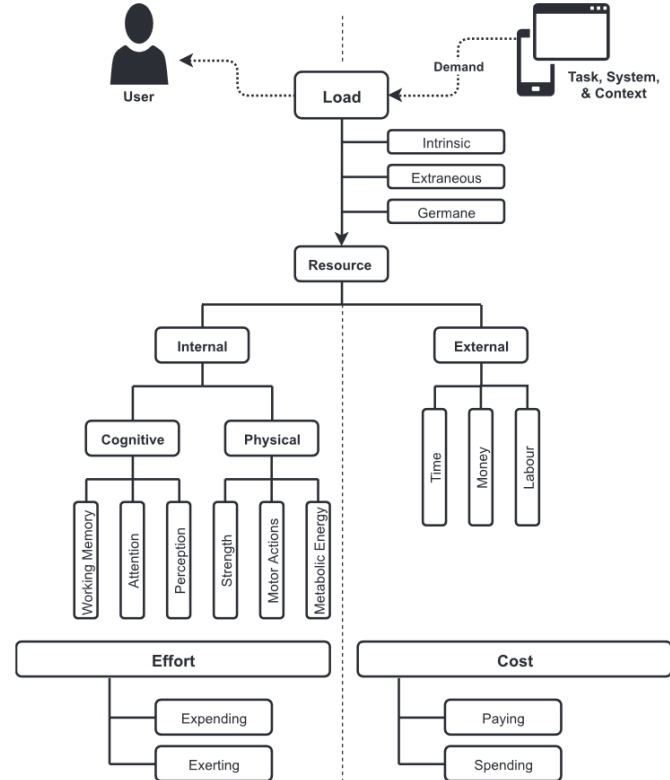(3) Estimating Cost of IR (discussion)

# Estimating Cost of IR
## Starting the Discussion

What do we mean by cost?

cost ➜ system (time, money, energy)

- ❑ training efficiency?
- ❑ inference efficiency?
- ❑ energy utilisation?



cost ➜ user (cost, effort, load) [McGregor et al. 2023]

- ❑ cognitive costs, fatigue, spend or conserve my resources to achieve goal?
- ❑ cognitive or physical effort, task complexity, total labour/time to achieve goal
- ❑ cognitive load, demands, properties of task that regulate exertion, overload