

Compression and Contraction

Adway Girish
Information Theory Lab

EPFL



September 9, 2024
IPG PhD Review

Outline

- 1 Prompt compression for black-box language models
- 2 Input-entropy-constrained capacity
- 3 Joint range of divergences
- 4 Closing remarks

Outline

- 1 Prompt compression for black-box language models
- 2 Input-entropy-constrained capacity
- 3 Joint range of divergences
- 4 Closing remarks

Prompt compression



Prompt compression

Prompt

It was the best of times, it was the worst of times, it was the age of wisdom, it was the age of foolishness, it was the epoch of belief, it was the epoch of incredulity, it was the season of light, it was the season of darkness, it was the spring of hope, it was the winter of despair.

x

LLM

Prompt compression

Prompt

It was the best of times, it was the worst of times, it was the age of wisdom, it was the age of foolishness, it was the epoch of belief, it was the epoch of incredulity, it was the season of light, it was the season of darkness, it was the spring of hope, it was the winter of despair.

x

LLM

q

Query

How were the times?

Prompt compression

Prompt

It was the best of times, it was the worst of times, it was the age of wisdom, it was the age of foolishness, it was the epoch of belief, it was the epoch of incredulity, it was the season of light, it was the season of darkness, it was the spring of hope, it was the winter of despair.

x

LLM

$$P_{\hat{y}} = \phi_{\text{LLM}}(x, q)$$

q

Query

How were the times?

Output

Best and worst.	(60%)
Contrasting.	(20%)
Mixed.	(10%)
Dualistic.	(5%)
⋮	

Prompt compression: query-agnostic

Prompt

It was the best of times, it was the worst of times, it was the age of wisdom, it was the age of foolishness, it was the epoch of belief, it was the epoch of incredulity, it was the season of light, it was the season of darkness, it was the spring of hope, it was the winter of despair.

Compressed prompt (query-agnostic)

best times worst, age wisdom foolish, epoch belief incredul, season light dark, hope despair.

x

comp

m

LLM

$P_{\hat{y}} = \phi_{\text{LLM}}(m, q)$

q

Query

How were the times?

Output

Best and worst.	(60%)
Contrasting.	(20%)
Mixed.	(10%)
Dualistic.	(5%)
⋮	

Prompt compression: query-aware

Prompt

It was the best of times, it was the worst of times, it was the age of wisdom, it was the age of foolishness, it was the epoch of belief, it was the epoch of incredulity, it was the season of light, it was the season of darkness, it was the spring of hope, it was the winter of despair.

x



m

Compressed prompt (query-aware)

best worst.



$P_{\hat{y}} = \phi_{\text{LLM}}(m, q)$

q

Query

How were the times?

Output

Best and worst.	(60%)
Contrasting.	(20%)
Mixed.	(10%)
Dualistic.	(5%)
⋮	

Prompt compression: rate-distortion formulation

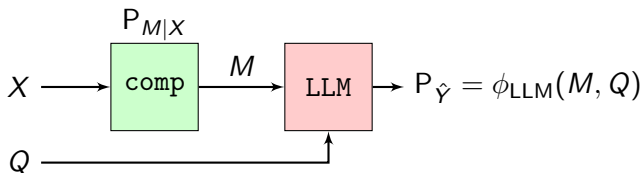
- $(X, Q, Y) \sim P_{XQY} = P_{XQ} P_{Y|XQ}$

$Y = \text{"true answer"}$

Prompt compression: rate-distortion formulation

- $(X, Q, Y) \sim P_{XQY} = P_{XQ} P_{Y|XQ}$

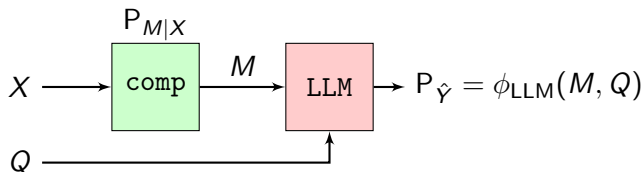
$Y = \text{"true answer"}$



Prompt compression: rate-distortion formulation

- $(X, Q, Y) \sim P_{XQY} = P_{XQ} P_{Y|XQ}$

$Y = \text{"true answer"}$

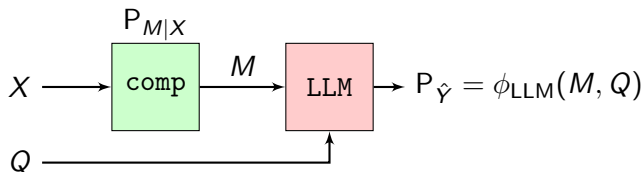


- Compression with side-information

Prompt compression: rate-distortion formulation

- $(X, Q, Y) \sim P_{XQY} = P_{XQ} P_{Y|XQ}$

$Y = \text{"true answer"}$

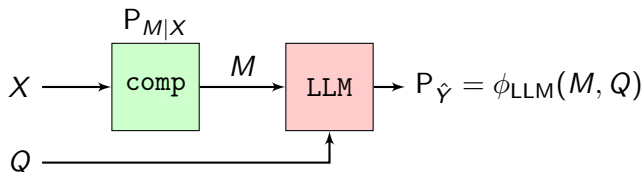


- Compression with side-information
for a **fixed decoder**, " $(m, q) \mapsto \phi_{\text{LLM}}(m, q)$ "

Prompt compression: rate-distortion formulation

- $(X, Q, Y) \sim P_{XQY} = P_{XQ} P_{Y|XQ}$

$Y = \text{"true answer"}$

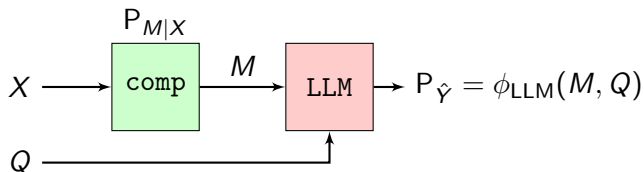


- Compression with side-information
for a **fixed decoder**, " $(m, q) \mapsto \phi_{\text{LLM}}(m, q)$ "
- Performance metrics:

Prompt compression: rate-distortion formulation

- $(X, Q, Y) \sim P_{XQY} = P_{XQ} P_{Y|XQ}$

$Y = \text{"true answer"}$



- Compression with side-information
for a **fixed decoder**, “ $(m, q) \mapsto \phi_{LLM}(m, q)$ ”
- Performance metrics:

$$\text{rate} = \mathbb{E} \left[\frac{\text{len}(M)}{\text{len}(X)} \right] \quad \text{distortion} = \mathbb{E} [d(Y, \phi_{LLM}(M, Q))]$$

Distortion-rate function

- $\text{rate} = \mathbb{E} \left[\frac{\text{len}(M)}{\text{len}(X)} \right] \quad \text{distortion} = \mathbb{E} [d(Y, \phi_{\text{LLM}}(M, Q))]$

Distortion-rate function

- $\text{rate} = \mathbb{E} \left[\frac{\text{len}(M)}{\text{len}(X)} \right] \quad \text{distortion} = \mathbb{E} [d(Y, \phi_{\text{LLM}}(M, Q))]$

- $$D^*(R) = \inf_{P_{M|X}} \mathbb{E} [d(Y, \phi_{\text{LLM}}(M, Q))]$$

s.t. $\mathbb{E} \left[\frac{\text{len}(M)}{\text{len}(X)} \right] \leq R$, and
 $P_{M|X}$ “is a compressor”

Distortion-rate function

- $\text{rate} = \mathbb{E} \left[\frac{\text{len}(M)}{\text{len}(X)} \right] \quad \text{distortion} = \mathbb{E} [d(Y, \phi_{\text{LLM}}(M, Q))]$

- $$D^*(R) = \inf_{P_{M|X}} \mathbb{E} [d(Y, \phi_{\text{LLM}}(M, Q))]$$

s.t. $\mathbb{E} \left[\frac{\text{len}(M)}{\text{len}(X)} \right] \leq R$, and
 $P_{M|X}$ “is a compressor”

- Linear program, but large dimension

Distortion-rate function

- $\text{rate} = \mathbb{E} \left[\frac{\text{len}(M)}{\text{len}(X)} \right]$ $\text{distortion} = \mathbb{E} [d(Y, \phi_{\text{LLM}}(M, Q))]$
- $$D^*(R) = \inf_{P_{M|X}} \mathbb{E} [d(Y, \phi_{\text{LLM}}(M, Q))]$$

s.t. $\mathbb{E} \left[\frac{\text{len}(M)}{\text{len}(X)} \right] \leq R$, and
 $P_{M|X}$ “is a compressor”
- Linear program, but large dimension $\approx 32,000^{10}$

Distortion-rate function

- $\text{rate} = \mathbb{E} \left[\frac{\text{len}(M)}{\text{len}(X)} \right] \quad \text{distortion} = \mathbb{E} [d(Y, \phi_{\text{LLM}}(M, Q))]$

- $$D^*(R) = \inf_{P_{M|X}} \mathbb{E} [d(Y, \phi_{\text{LLM}}(M, Q))]$$

s.t. $\mathbb{E} \left[\frac{\text{len}(M)}{\text{len}(X)} \right] \leq R$, and
 $P_{M|X}$ “is a compressor”

- Linear program, but large dimension $\approx 32,000^{10}$

- Dual:

$$D^*(R) = \sup_{\lambda \geq 0} \left\{ -\lambda R + \sum_{x \in \mathcal{X}} \min_{m \in \mathcal{M}_x} \left[\mathbf{D}_{x,m} + \lambda \mathbf{R}_{x,m} \right] \right\}$$

Distortion-rate function

- $\text{rate} = \mathbb{E} \left[\frac{\text{len}(M)}{\text{len}(X)} \right] \quad \text{distortion} = \mathbb{E} [d(Y, \phi_{\text{LLM}}(M, Q))]$

- $$D^*(R) = \inf_{P_{M|X}} \mathbb{E} [d(Y, \phi_{\text{LLM}}(M, Q))]$$

s.t. $\mathbb{E} \left[\frac{\text{len}(M)}{\text{len}(X)} \right] \leq R$, and
 $P_{M|X}$ “is a compressor”

- Linear program, but large dimension $\approx 32,000^{10}$

- Dual:
$$D^*(R) = \sup_{\lambda \geq 0} \left\{ -\lambda R + \sum_{x \in \mathcal{X}} \min_{m \in \mathcal{M}_x} \left[\textcolor{red}{D}_{x,m} + \lambda \textcolor{red}{R}_{x,m} \right] \right\}$$

all possible “compressions” of x

Distortion-rate function

- $\text{rate} = \mathbb{E} \left[\frac{\text{len}(M)}{\text{len}(X)} \right] \quad \text{distortion} = \mathbb{E} [d(Y, \phi_{\text{LLM}}(M, Q))]$

- $$D^*(R) = \inf_{P_{M|X}} \mathbb{E} [d(Y, \phi_{\text{LLM}}(M, Q))]$$

s.t. $\mathbb{E} \left[\frac{\text{len}(M)}{\text{len}(X)} \right] \leq R$, and
 $P_{M|X}$ “is a compressor”

- Linear program, but large dimension $\approx 32,000^{10}$

- Dual:
$$D^*(R) = \sup_{\lambda \geq 0} \left\{ -\lambda R + \sum_{x \in \mathcal{X}} \min_{m \in \mathcal{M}_x} \left[\underset{\substack{\uparrow \\ \text{“normalized” distortion, rate} \\ \text{on compressing } x \mapsto m}}}{D_{x,m}} + \lambda \underset{\substack{\uparrow \\ \text{rate}}}{R_{x,m}} \right] \right\}$$

all possible “compressions” of x

Distortion-rate function: geometric solution via dual

- Dual:
$$D^*(R) = \sup_{\lambda \geq 0} \left\{ -\lambda R + \sum_{x \in \mathcal{X}} \min_{m \in \mathcal{M}_x} [D_{x,m} + \lambda R_{x,m}] \right\}$$

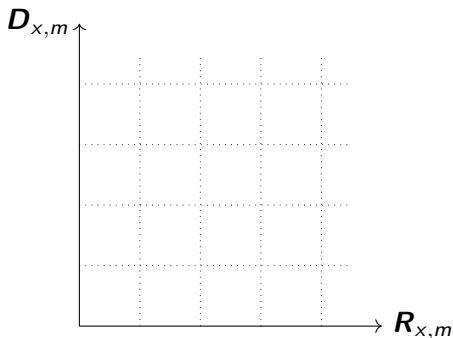
Distortion-rate function: geometric solution via dual

- Dual:
$$D^*(R) = \sup_{\lambda \geq 0} \left\{ -\lambda R + \sum_{x \in \mathcal{X}} \underbrace{\min_{m \in \mathcal{M}_x} [D_{x,m} + \lambda R_{x,m}]}_{\text{for fixed } (\lambda, x)} \right\}$$

Distortion-rate function: geometric solution via dual

- Dual:
$$D^*(R) = \sup_{\lambda \geq 0} \left\{ -\lambda R + \sum_{x \in \mathcal{X}} \underbrace{\min_{m \in \mathcal{M}_x} [D_{x,m} + \lambda R_{x,m}]}_{\text{for fixed } (\lambda, x)} \right\}$$

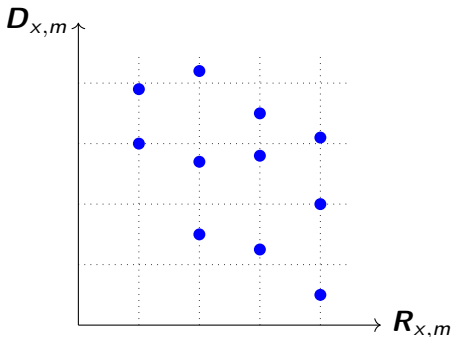
- Fix $\lambda \geq 0, x \in \mathcal{X}$



Distortion-rate function: geometric solution via dual

- Dual:
$$D^*(R) = \sup_{\lambda \geq 0} \left\{ -\lambda R + \sum_{x \in \mathcal{X}} \underbrace{\min_{m \in \mathcal{M}_x} [D_{x,m} + \lambda R_{x,m}]}_{\text{for fixed } (\lambda, x)} \right\}$$

- Fix $\lambda \geq 0$, $x \in \mathcal{X}$

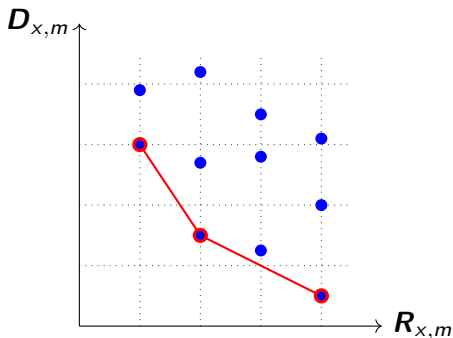


- Relevant points: 32,000¹⁰

Distortion-rate function: geometric solution via dual

- Dual:
$$D^*(R) = \sup_{\lambda \geq 0} \left\{ -\lambda R + \sum_{x \in \mathcal{X}} \underbrace{\min_{m \in \mathcal{M}_x} [D_{x,m} + \lambda R_{x,m}]}_{\text{for fixed } (\lambda, x)} \right\}$$

- Fix $\lambda \geq 0$, $x \in \mathcal{X}$

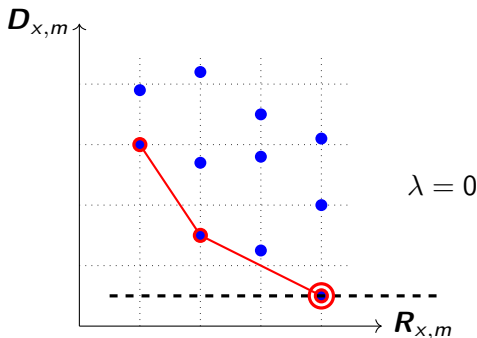


- Relevant points: 32,000¹⁰

Distortion-rate function: geometric solution via dual

- Dual:
$$D^*(R) = \sup_{\lambda \geq 0} \left\{ -\lambda R + \sum_{x \in \mathcal{X}} \underbrace{\min_{m \in \mathcal{M}_x} [D_{x,m} + \lambda R_{x,m}]}_{\text{for fixed } (\lambda, x)} \right\}$$

- Fix $\lambda \geq 0$, $x \in \mathcal{X}$

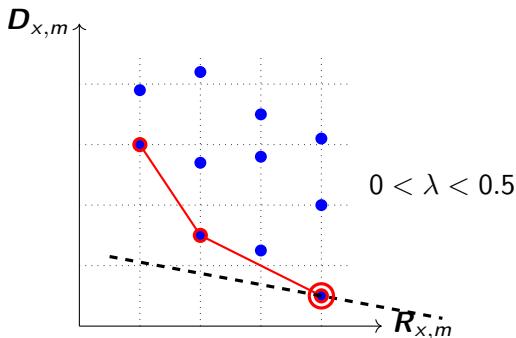


- Relevant points: 32,000¹⁰

Distortion-rate function: geometric solution via dual

- Dual:
$$D^*(R) = \sup_{\lambda \geq 0} \left\{ -\lambda R + \sum_{x \in \mathcal{X}} \underbrace{\min_{m \in \mathcal{M}_x} [D_{x,m} + \lambda R_{x,m}]}_{\text{for fixed } (\lambda, x)} \right\}$$

- Fix $\lambda \geq 0$, $x \in \mathcal{X}$

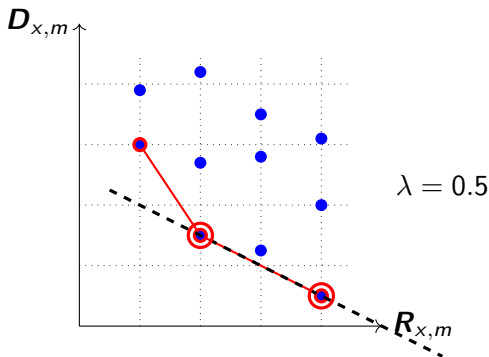


- Relevant points: 32,000¹⁰

Distortion-rate function: geometric solution via dual

- Dual:
$$D^*(R) = \sup_{\lambda \geq 0} \left\{ -\lambda R + \sum_{x \in \mathcal{X}} \underbrace{\min_{m \in \mathcal{M}_x} [D_{x,m} + \lambda R_{x,m}]}_{\text{for fixed } (\lambda, x)} \right\}$$

- Fix $\lambda \geq 0, x \in \mathcal{X}$

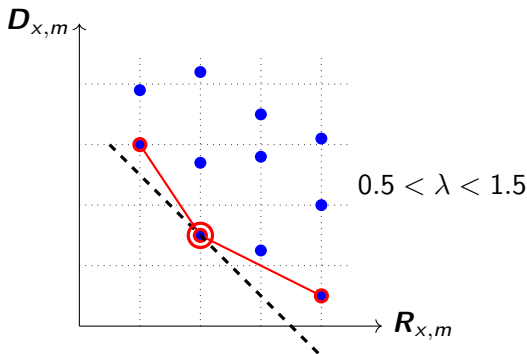


- Relevant points: 32,000¹⁰

Distortion-rate function: geometric solution via dual

- Dual:
$$D^*(R) = \sup_{\lambda \geq 0} \left\{ -\lambda R + \sum_{x \in \mathcal{X}} \underbrace{\min_{m \in \mathcal{M}_x} [D_{x,m} + \lambda R_{x,m}]}_{\text{for fixed } (\lambda, x)} \right\}$$

- Fix $\lambda \geq 0$, $x \in \mathcal{X}$

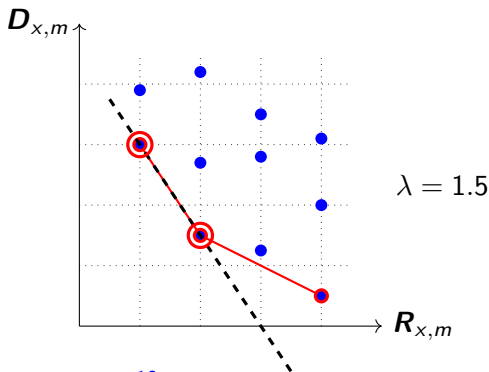


- Relevant points: 32,000¹⁰

Distortion-rate function: geometric solution via dual

- Dual:
$$D^*(R) = \sup_{\lambda \geq 0} \left\{ -\lambda R + \sum_{x \in \mathcal{X}} \underbrace{\min_{m \in \mathcal{M}_x} [D_{x,m} + \lambda R_{x,m}]}_{\text{for fixed } (\lambda, x)} \right\}$$

- Fix $\lambda \geq 0, x \in \mathcal{X}$

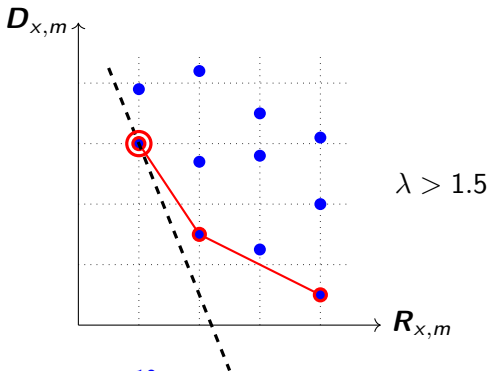


- Relevant points: 32,000¹⁰

Distortion-rate function: geometric solution via dual

- Dual:
$$D^*(R) = \sup_{\lambda \geq 0} \left\{ -\lambda R + \sum_{x \in \mathcal{X}} \underbrace{\min_{m \in \mathcal{M}_x} [D_{x,m} + \lambda R_{x,m}]}_{\text{for fixed } (\lambda, x)} \right\}$$

- Fix $\lambda \geq 0$, $x \in \mathcal{X}$

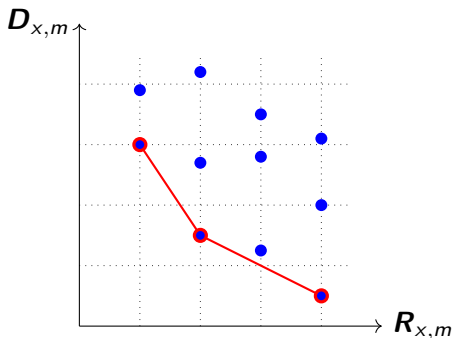


- Relevant points: $32,000^{10}$

Distortion-rate function: geometric solution via dual

- Dual:
$$D^*(R) = \sup_{\lambda \geq 0} \left\{ -\lambda R + \sum_{x \in \mathcal{X}} \underbrace{\min_{m \in \mathcal{M}_x} [D_{x,m} + \lambda R_{x,m}]}_{\text{for fixed } (\lambda, x)} \right\}$$

- Fix $\lambda \geq 0$, $x \in \mathcal{X}$

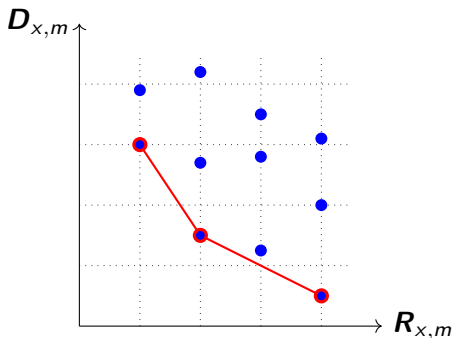


- Relevant points: $32,000^{10} \rightarrow 10$

Distortion-rate function: geometric solution via dual

- Dual:
$$D^*(R) = \sup_{\lambda \geq 0} \left\{ -\lambda R + \sum_{x \in \mathcal{X}} \underbrace{\min_{m \in \mathcal{M}_x} [D_{x,m} + \lambda R_{x,m}]}_{\text{for fixed } (\lambda, x)} \right\}$$

- Fix $\lambda \geq 0$, $x \in \mathcal{X}$

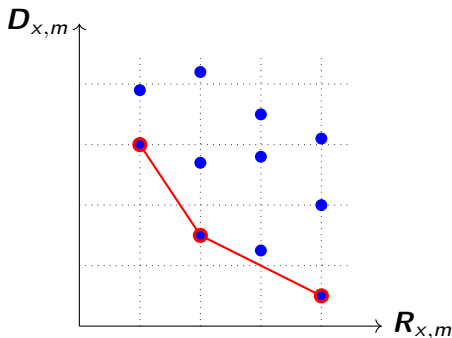


- Relevant points: $32,000^{10} \rightarrow 10$,
only finitely many λ

Distortion-rate function: geometric solution via dual

- Dual:
$$D^*(R) = \sup_{\lambda \geq 0} \left\{ -\lambda R + \sum_{x \in \mathcal{X}} \underbrace{\min_{m \in \mathcal{M}_x} [D_{x,m} + \lambda R_{x,m}]}_{\text{for fixed } (\lambda, x)} \right\}$$

- Fix $\lambda \geq 0$, $x \in \mathcal{X}$

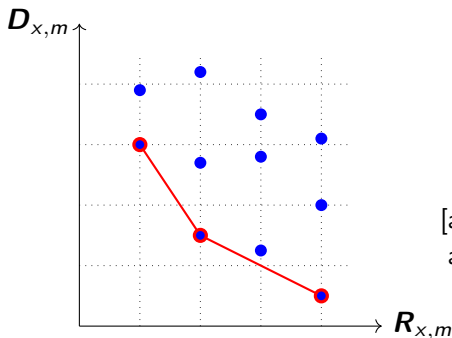


- Relevant points: $32,000^{10} \rightarrow 10$, $(2^{10} \rightarrow 10)$
only finitely many λ

Distortion-rate function: geometric solution via dual

- Dual:
$$D^*(R) = \sup_{\lambda \geq 0} \left\{ -\lambda R + \sum_{x \in \mathcal{X}} \underbrace{\min_{m \in \mathcal{M}_x} [D_{x,m} + \lambda R_{x,m}]}_{\text{for fixed } (\lambda, x)} \right\}$$

- Fix $\lambda \geq 0$, $x \in \mathcal{X}$

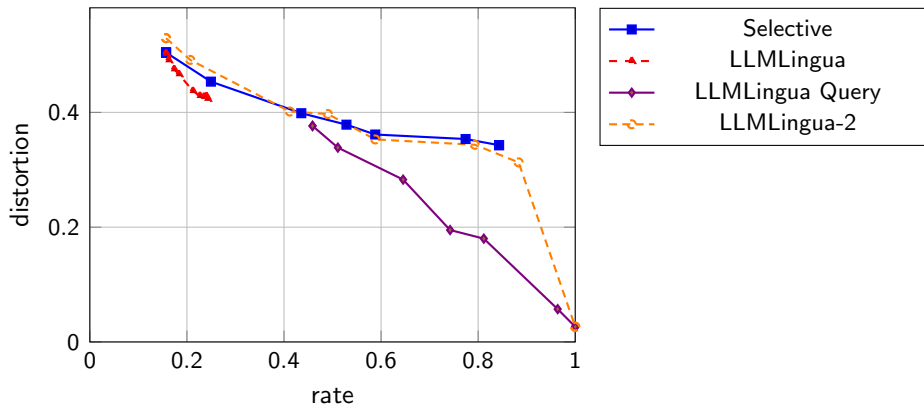


[apple \mapsto app, ale, pe;
apple $\not\mapsto$ pale, red, lp]

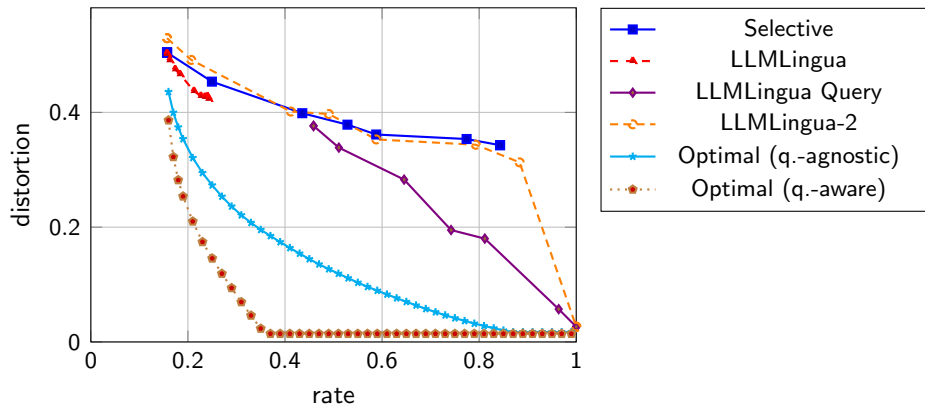
- Relevant points: $32,000^{10} \rightarrow 10$,
only finitely many λ

$(2^{10} \rightarrow 10)$

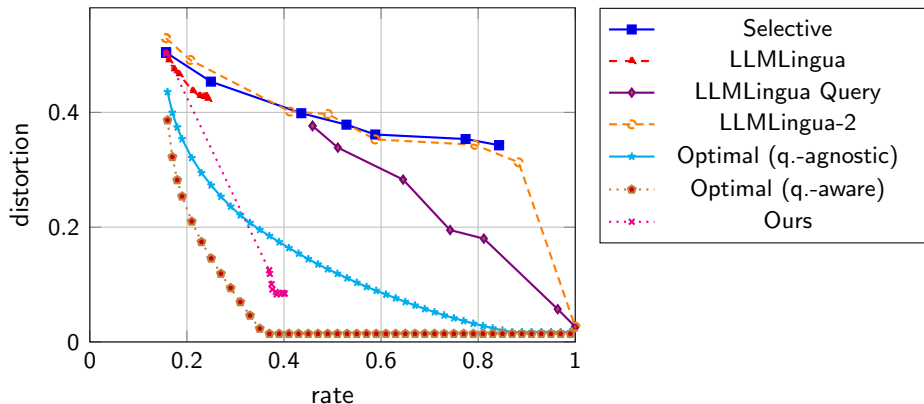
Experimental results



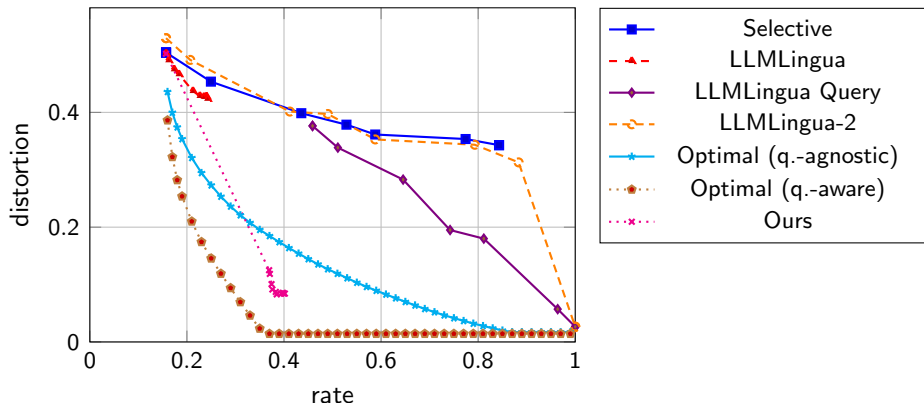
Experimental results



Experimental results



Experimental results



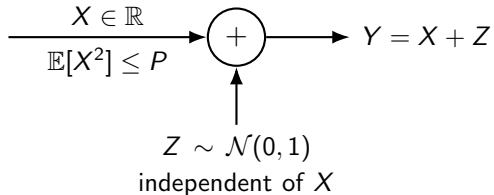
A.G.*, A.Nagle*, M.Bondaschi, M.Gastpar, A.V.Makkuva, H.Kim, “Fundamental Limits of Prompt Compression: A Rate-Distortion Framework for Black-Box Language Models.”
— ICML 2024 Workshop on Theoretical Foundations of Foundation Models [Oral]
— under review at NeurIPS 2024

Segue to a contraction problem

- Optimization 101...

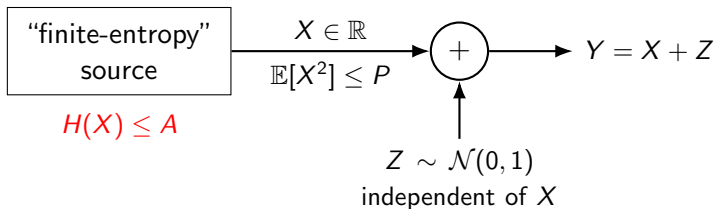
Segue to a contraction problem

- Optimization 101... thanks to a different problem:



Segue to a contraction problem

- Optimization 101... thanks to a different problem:

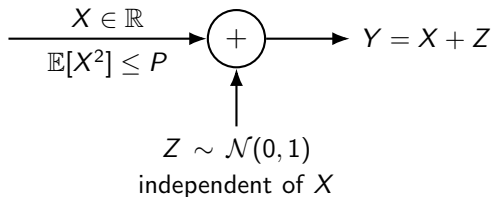


Outline

- 1 Prompt compression for black-box language models
- 2 Input-entropy-constrained capacity**
- 3 Joint range of divergences
- 4 Closing remarks

Input-entropy-constrained channel capacity

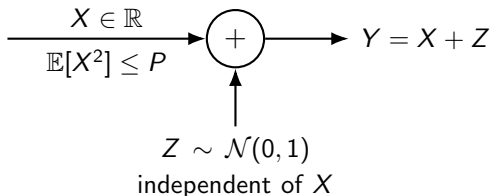
- A contraction problem in communication



$$C_H(A, P) = \sup_{\substack{P_X: \mathbb{E}[X^2] \leq P \\ H(X) \leq A}} I(X; Y)$$

Input-entropy-constrained channel capacity

- A contraction problem in communication

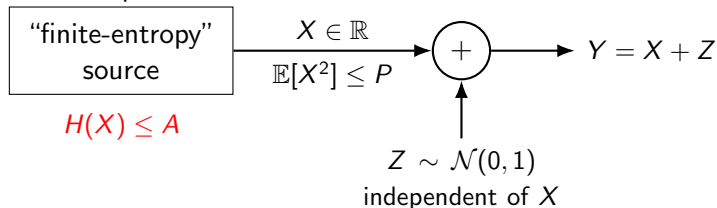


$$C_H(A, P) = \sup_{\substack{P_X: \mathbb{E}[X^2] \leq P \\ H(X) \leq A}} I(X; Y)$$

- Cardinality bounds?

Input-entropy-constrained channel capacity

- A contraction problem in communication

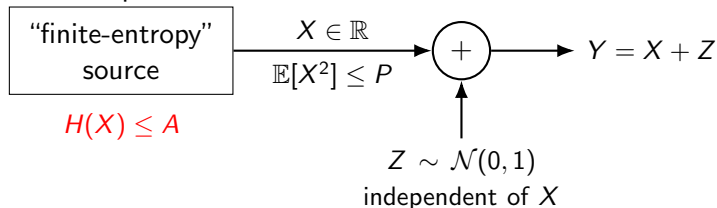


$$C_H(A, P) = \sup_{\substack{P_X: \\ \mathbb{E}[X^2] \leq P \\ H(X) \leq A}} I(X; Y)$$

- Cardinality bounds? Finite support?

Input-entropy-constrained channel capacity

- A contraction problem in communication



$$C_H(A, P) = \sup_{\substack{P_X: \\ \mathbb{E}[X^2] \leq P \\ H(X) \leq A}} I(X; Y)$$

- Cardinality bounds? Finite support?
- A nontrivial upper bound better than

$$F_I(A, P) = \sup_{\substack{P_{WX}: \\ \mathbb{E}[X^2] \leq P \\ I(W; X) \leq A}} I(W; Y) \quad ?$$

Aside on data processing inequalities

Fix $P_{Y|X}$

Aside on data processing inequalities

Fix $P_{Y|X}$

- DPI: for any P_{WX} , $I(W; Y) \leq I(W; X)$

Aside on data processing inequalities

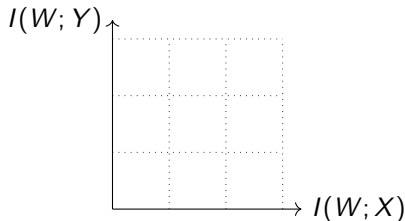
Fix $P_{Y|X}$

- DPI: for any P_{WX} , $I(W; Y) \leq I(W; X)$
- Data processing *function*: $F_I(t) = \sup_{P_{WX}: I(W; X) \leq t} I(W; Y)$

Aside on data processing inequalities

Fix $P_{Y|X}$

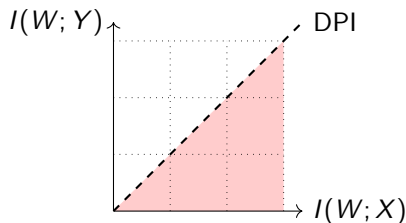
- DPI: for any P_{WX} , $I(W; Y) \leq I(W; X)$
- Data processing *function*: $F_I(t) = \sup_{P_{WX}: I(W; X) \leq t} I(W; Y)$



Aside on data processing inequalities

Fix $P_{Y|X}$

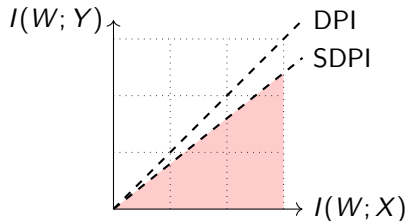
- DPI: for any P_{WX} , $I(W; Y) \leq I(W; X)$
- Data processing *function*: $F_I(t) = \sup_{P_{WX}: I(W; X) \leq t} I(W; Y)$



Aside on data processing inequalities

Fix $P_{Y|X}$

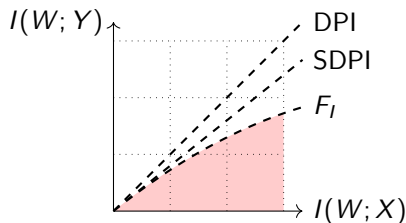
- DPI: for any P_{WX} , $I(W; Y) \leq I(W; X)$
- Data processing *function*: $F_I(t) = \sup_{P_{WX}: I(W; X) \leq t} I(W; Y)$



Aside on data processing inequalities

Fix $P_{Y|X}$

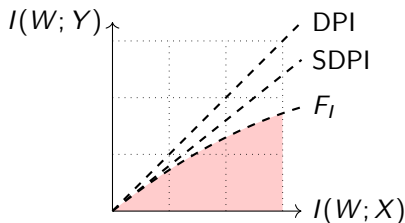
- DPI: for any P_{WX} , $I(W; Y) \leq I(W; X)$
- Data processing *function*: $F_I(t) = \sup_{P_{WX}: I(W; X) \leq t} I(W; Y)$



Aside on data processing inequalities

Fix $P_{Y|X}$

- DPI: for any P_{WX} , $I(W; Y) \leq I(W; X)$
- Data processing function: $F_I(t) = \sup_{P_{WX}: I(W; X) \leq t} I(W; Y)$

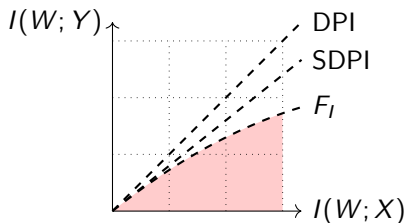


- Also DPI: for any P_X, Q_X , $D_f(Q_Y || P_Y) \leq D_f(Q_X || P_X)$

Aside on data processing inequalities

Fix $P_{Y|X}$

- DPI: for any P_{WX} , $I(W; Y) \leq I(W; X)$
- Data processing function: $F_I(t) = \sup_{P_{WX}: I(W; X) \leq t} I(W; Y)$

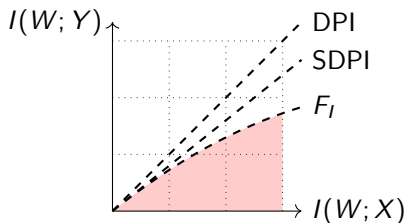


- Also DPI: for any P_X, Q_X , $D_f(Q_Y \| P_Y) \leq D_f(Q_X \| P_X)$
 $\nearrow P_Y = P_X \circ P_{Y|X}$

Aside on data processing inequalities

Fix $P_{Y|X}$

- DPI: for any P_{WX} , $I(W; Y) \leq I(W; X)$
- Data processing *function*: $F_I(t) = \sup_{P_{WX}: I(W; X) \leq t} I(W; Y)$

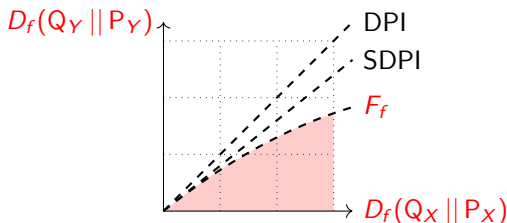


- Also DPI: for any P_X, Q_X , $D_f(Q_Y \parallel P_Y) \leq D_f(Q_X \parallel P_X)$
- Natural analogue: $F_f(t) = \sup_{P_X, Q_X: D_f(Q_X \parallel P_X) \leq t} D_f(Q_Y \parallel P_Y)$

Aside on data processing inequalities

Fix $P_{Y|X}$

- DPI: for any P_{WX} , $I(W; Y) \leq I(W; X)$
- Data processing *function*: $F_I(t) = \sup_{P_{WX}: I(W; X) \leq t} I(W; Y)$



- Also DPI: for any P_X, Q_X , $D_f(Q_Y || P_Y) \leq D_f(Q_X || P_X)$
- Natural analogue: $F_f(t) = \sup_{P_X, Q_X: D_f(Q_X || P_X) \leq t} D_f(Q_Y || P_Y)$

Outline

- 1 Prompt compression for black-box language models
- 2 Input-entropy-constrained capacity
- 3 Joint range of divergences**
- 4 Closing remarks

Joint range of input and output divergences

Fix $P_{Y|X}$ and f

Joint range of input and output divergences

Fix $P_{Y|X}$ and f

- Define $F_f(t) = \sup_{P_X, Q_X: D_f(Q_X || P_X) \leq t} D_f(Q_Y || P_Y)$

Joint range of input and output divergences

Fix $P_{Y|X}$ and f

- Define $F_f(t) = \sup_{P_X, Q_X: D_f(Q_X || P_X) \leq t} D_f(Q_Y || P_Y)$
- Upper boundary of $\mathcal{D}_f = \bigcup_{P_X, Q_X} \{ (D_f(Q_X || P_X), D_f(Q_Y || P_Y)) \}$

Joint range of input and output divergences

Fix $P_{Y|X}$ and f

- Define $F_f(t) = \sup_{P_X, Q_X: D_f(Q_X || P_X) \leq t} D_f(Q_Y || P_Y)$
- Upper boundary of $\mathcal{D}_f = \bigcup_{P_X, Q_X} \{ (D_f(Q_X || P_X), D_f(Q_Y || P_Y)) \}$
- Conjecture: \mathcal{D}_f is convex

?

Joint range of input and output divergences

Fix $P_{Y|X}$ and f

- Define $F_f(t) = \sup_{P_X, Q_X: D_f(Q_X || P_X) \leq t} D_f(Q_Y || P_Y)$
- Upper boundary of $\mathcal{D}_f = \bigcup_{P_X, Q_X} \{ (D_f(Q_X || P_X), D_f(Q_Y || P_Y)) \}$
- Conjecture: \mathcal{D}_f is convex ($\implies F_f$ is concave) ?

Joint range of input and output divergences

Fix $P_{Y|X}$ and f

- Define $F_f(t) = \sup_{P_X, Q_X: D_f(Q_X || P_X) \leq t} D_f(Q_Y || P_Y)$
- Upper boundary of $\mathcal{D}_f = \bigcup_{P_X, Q_X} \{ (D_f(Q_X || P_X), D_f(Q_Y || P_Y)) \}$
- Conjecture: \mathcal{D}_f is convex ($\implies F_f$ is concave) ?
- Facts:

Joint range of input and output divergences

Fix $P_{Y|X}$ and f

- Define $F_f(t) = \sup_{P_X, Q_X: D_f(Q_X || P_X) \leq t} D_f(Q_Y || P_Y)$
- Upper boundary of $\mathcal{D}_f = \bigcup_{P_X, Q_X} \{ (D_f(Q_X || P_X), D_f(Q_Y || P_Y)) \}$
- Conjecture: \mathcal{D}_f is convex ($\implies F_f$ is concave) ?
- Facts:
 - F_f is NOT necessarily concave (counter-example: $P_{Y|X} = \text{BEC}^3$)

Joint range of input and output divergences

Fix $P_{Y|X}$ and f

- Define $F_f(t) = \sup_{P_X, Q_X: D_f(Q_X || P_X) \leq t} D_f(Q_Y || P_Y)$
- Upper boundary of $\mathcal{D}_f = \bigcup_{P_X, Q_X} \{ (D_f(Q_X || P_X), D_f(Q_Y || P_Y)) \}$
- Conjecture: \mathcal{D}_f is convex ($\implies F_f$ is concave) ?
- Facts:
 - F_I is NOT necessarily concave (counter-example: $P_{Y|X} = \text{BEC}^3$)
 - Fix P_X , define $\tilde{F}_I(t, P_X) = \sup_{P_{W|X}: I(W; X) \leq t} I(W; Y)$ and

Joint range of input and output divergences

Fix $P_{Y|X}$ and f

- Define $F_f(t) = \sup_{P_X, Q_X: D_f(Q_X || P_X) \leq t} D_f(Q_Y || P_Y)$
- Upper boundary of $\mathcal{D}_f = \bigcup_{P_X, Q_X} \{ (D_f(Q_X || P_X), D_f(Q_Y || P_Y)) \}$
- Conjecture: \mathcal{D}_f is convex ($\implies F_f$ is concave) ?
- Facts:
 - F_I is NOT necessarily concave (counter-example: $P_{Y|X} = \text{BEC}^3$)
 - Fix P_X , define $\tilde{F}_I(t, P_X) = \sup_{P_{W|X}: I(W; X) \leq t} I(W; Y)$ and
$$\tilde{F}_f(t, P_X) = \sup_{Q_X: D_f(Q_X || P_X) \leq t} D_f(Q_Y || P_Y),$$

Joint range of input and output divergences

Fix $P_{Y|X}$ and f

- Define $F_f(t) = \sup_{P_X, Q_X: D_f(Q_X || P_X) \leq t} D_f(Q_Y || P_Y)$
- Upper boundary of $\mathcal{D}_f = \bigcup_{P_X, Q_X} \{ (D_f(Q_X || P_X), D_f(Q_Y || P_Y)) \}$
- Conjecture: \mathcal{D}_f is convex ($\implies F_f$ is concave) ?
- Facts:
 - F_I is NOT necessarily concave (counter-example: $P_{Y|X} = \text{BEC}^3$)
 - Fix P_X , define $\tilde{F}_I(t, P_X) = \sup_{P_{W|X}: I(W; X) \leq t} I(W; Y)$ and
$$\tilde{F}_f(t, P_X) = \sup_{Q_X: D_f(Q_X || P_X) \leq t} D_f(Q_Y || P_Y),$$
then $\tilde{F}_I(\cdot, P_X)$ is concave

Joint range of input and output divergences

Fix $P_{Y|X}$ and f

- Define $F_f(t) = \sup_{P_X, Q_X: D_f(Q_X || P_X) \leq t} D_f(Q_Y || P_Y)$
- Upper boundary of $\mathcal{D}_f = \bigcup_{P_X, Q_X} \{ (D_f(Q_X || P_X), D_f(Q_Y || P_Y)) \}$
- Conjecture: \mathcal{D}_f is convex ($\implies F_f$ is concave) ?
- Facts:
 - F_I is NOT necessarily concave (counter-example: $P_{Y|X} = \text{BEC}^3$)
 - Fix P_X , define $\tilde{F}_I(t, P_X) = \sup_{P_{W|X}: I(W; X) \leq t} I(W; Y)$ and
$$\tilde{F}_f(t, P_X) = \sup_{Q_X: D_f(Q_X || P_X) \leq t} D_f(Q_Y || P_Y),$$
then $\tilde{F}_I(\cdot, P_X)$ is concave; $\implies \tilde{F}_f(\cdot, P_X)$ is concave

Joint range of input and output divergences

Fix $P_{Y|X}$ and f

- Define $F_f(t) = \sup_{P_X, Q_X: D_f(Q_X || P_X) \leq t} D_f(Q_Y || P_Y)$
- Upper boundary of $\mathcal{D}_f = \bigcup_{P_X, Q_X} \{ (D_f(Q_X || P_X), D_f(Q_Y || P_Y)) \}$
- Conjecture: \mathcal{D}_f is convex ($\implies F_f$ is concave) ?
- Facts:
 - F_f is NOT necessarily concave (counter-example: $P_{Y|X} = \text{BEC}^3$)
 - Fix P_X , define $\tilde{F}_f(t, P_X) = \sup_{P_{W|X}: I(W; Y) \leq t} I(W; Y)$ and
$$\tilde{F}_f(t, P_X) = \sup_{Q_X: D_f(Q_X || P_X) \leq t} D_f(Q_Y || P_Y),$$
then $\tilde{F}_f(\cdot, P_X)$ is concave; $\implies \tilde{F}_f(\cdot, P_X)$ is concave
 - For any f, g , $\bigcup_{P_X, Q_X} \{ (D_f(Q_X || P_X), D_g(Q_X || P_X)) \}$ is convex

Outline

- 1 Prompt compression for black-box language models
- 2 Input-entropy-constrained capacity
- 3 Joint range of divergences
- 4 Closing remarks

In closing...

In closing...

- Three problems ($1\times$ compression, $2\times$ contraction):

In closing...

- Three problems ($1\times$ compression, $2\times$ contraction):
 - Prompt compression for LLMs

In closing...

- Three problems ($1\times$ compression, $2\times$ contraction):
 - Prompt compression for LLMs
 - Entropy-constrained capacity

In closing...

- Three problems ($1\times$ compression, $2\times$ contraction):
 - Prompt compression for LLMs
 - Entropy-constrained capacity
 - Joint range of divergences

In closing...

- Three problems ($1\times$ compression, $2\times$ contraction):
 - Prompt compression for LLMs
 - Entropy-constrained capacity
 - Joint range of divergences
- Two more:

In closing...

- Three problems ($1\times$ compression, $2\times$ contraction):
 - Prompt compression for LLMs
 - Entropy-constrained capacity
 - Joint range of divergences
- Two more:
 - Guesswork

In closing...

- Three problems ($1\times$ compression, $2\times$ contraction):
 - Prompt compression for LLMs
 - Entropy-constrained capacity
 - Joint range of divergences
- Two more:
 - Guesswork
 - Distributed hypothesis testing

In closing...

- Three problems ($1\times$ compression, $2\times$ contraction):
 - Prompt compression for LLMs
 - Entropy-constrained capacity
 - Joint range of divergences
- Two more (method of types + optimization):
 - Guesswork
 - Distributed hypothesis testing

In closing...

- Three problems ($1\times$ compression, $2\times$ contraction):
 - Prompt compression for LLMs
 - Entropy-constrained capacity
 - Joint range of divergences
- Two more (method of types + optimization):
 - Guesswork
 - Distributed hypothesis testing \rightarrow compression + contraction

In closing...

- Three problems ($1\times$ compression, $2\times$ contraction):
 - Prompt compression for LLMs
 - Entropy-constrained capacity
 - Joint range of divergences
- Two more (method of types + optimization):
 - Guesswork
 - Distributed hypothesis testing \rightarrow compression + contraction
- All thoughts welcome

In closing...

- Three problems ($1\times$ compression, $2\times$ contraction):
 - Prompt compression for LLMs
 - Entropy-constrained capacity
 - Joint range of divergences
- Two more (method of types + optimization):
 - Guesswork
 - Distributed hypothesis testing \longrightarrow compression + contraction
- All thoughts welcome

Thank you!