# Compression and Contraction

Adway Girish
Information Theory Lab

**EPFL**

EPFL Information Processing Group

September 9, 2024
IPG PhD Review

# Outline

# Outline

# 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 ⟶ $P_{\hat{Y}} = \phi_{\text{LLM}}(x, q)$

**Output**

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

⋮

$q$

**Query**

How were the times?

# 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.

$x \longrightarrow$ `comp` $\xrightarrow{m}$ `LLM` $\longrightarrow P_{\hat{Y}} = \phi_{\mathsf{LLM}}(m, q)$

**Output**

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

$\vdots$

$q$ ——

**Query**

How were the times?

# 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.

**Compressed prompt** (query-aware)

best worst.

$x \longrightarrow$ comp $\xrightarrow{\quad m \quad}$ LLM $\longrightarrow$ $\mathsf{P}_{\hat{Y}} = \phi_{\mathsf{LLM}}(m, q)$

**Output**

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

$q$

**Query**

How were the times?

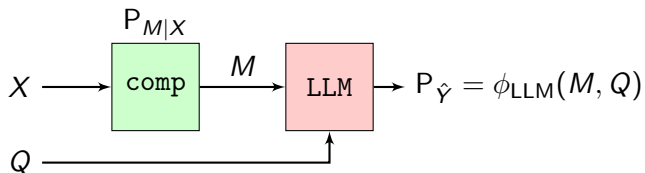# Prompt compression: rate-distortion formulation

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

# Prompt compression: rate-distortion formulation

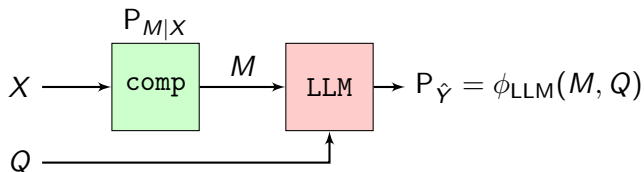- $(X, Q, Y) \sim P_{XQY} = P_{XQ} P_{Y|XQ}$ $\qquad\qquad$ $Y = $ "true answer"



$$X \longrightarrow \boxed{\text{comp}} \xrightarrow{M} \boxed{\text{LLM}} \longrightarrow P_{\hat{Y}} = \phi_{\mathsf{LLM}}(M, Q)$$

with $P_{M|X}$ above the comp box and $Q$ feeding into the LLM box.

# Prompt compression: rate-distortion formulation

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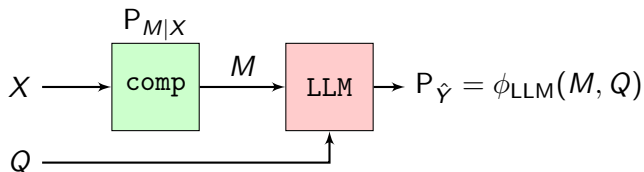


- Compression with side-information

# Prompt compression: rate-distortion formulation

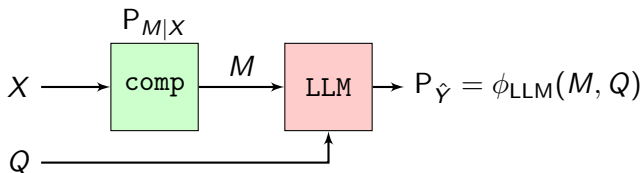- $(X, Q, Y) \sim P_{XQY} = P_{XQ} \, P_{Y|XQ}$ $\qquad\qquad$ $Y =$ "true answer"



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

# Prompt compression: rate-distortion formulation

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

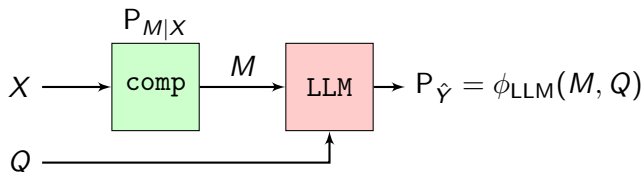

- Compression with side-information
  for a fixed decoder, "$(m, q) \mapsto \phi_{\mathsf{LLM}}(m, q)$"

- Performance metrics:

# Prompt compression: rate-distortion formulation

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



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

- Performance metrics:

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

# Distortion-rate function

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

# Distortion-rate function

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

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

  $\text{s.t.}$ $\mathbb{E}\left[\frac{\text{len}(M)}{\text{len}(X)}\right] \leq R$, and

  $P_{M|X}$ "is a compressor"

# Distortion-rate function

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

- $$D^*(R) = \inf_{\mathsf{P}_{M|X}} \quad \mathbb{E}\left[\text{d}\big(Y, \phi_{\text{LLM}}(M, Q)\big)\right]$$
  $$\text{s.t.} \quad \mathbb{E}\left[\frac{\text{len}(M)}{\text{len}(X)}\right] \leq R, \text{ and}$$
  $$\mathsf{P}_{M|X} \text{ "is a compressor"}$$

- Linear program, but large dimension

# Distortion-rate function

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

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

    $$\text{s.t.} \quad \mathbb{E}\left[\frac{\text{len}(M)}{\text{len}(X)}\right] \leq R, \text{ and}$$

    $$P_{M|X} \text{ "is a compressor"}$$

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

# Distortion-rate function

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

- $$D^*(R) = \inf_{\mathsf{P}_{M|X}} \quad \mathbb{E}\left[\text{d}\big(Y, \phi_{\text{LLM}}(M, Q)\big)\right]$$
  $$\text{s.t.} \quad \mathbb{E}\left[\frac{\text{len}(M)}{\text{len}(X)}\right] \leq R, \text{ and}$$
  $$\mathsf{P}_{M|X} \text{ "is a compressor"}$$

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

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

# Distortion-rate function

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

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

$$\text{s.t.} \quad \mathbb{E}\left[\frac{\text{len}(M)}{\text{len}(X)}\right] \leq R, \text{ and}$$

$$P_{M|X} \text{ "is a compressor"}$$

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

- Dual:

all possible "compressions" of $x$

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

# Distortion-rate function

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

- $$D^*(R) = \inf_{\mathsf{P}_{M|X}} \quad \mathbb{E}\left[\mathsf{d}(Y, \phi_{\text{LLM}}(M, Q))\right]$$
  $$\text{s.t.} \quad \mathbb{E}\left[\frac{\text{len}(M)}{\text{len}(X)}\right] \leq R, \text{ and}$$
  $$\mathsf{P}_{M|X} \text{ "is a compressor"}$$

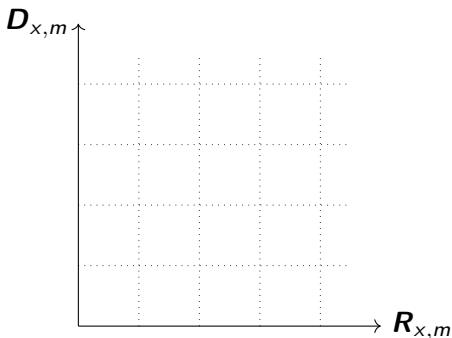- Linear program, but large dimension $\approx 32{,}000^{10}$

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

  all possible "compressions" of $x$

  "normalized" distortion, rate
  on compressing $x \mapsto m$
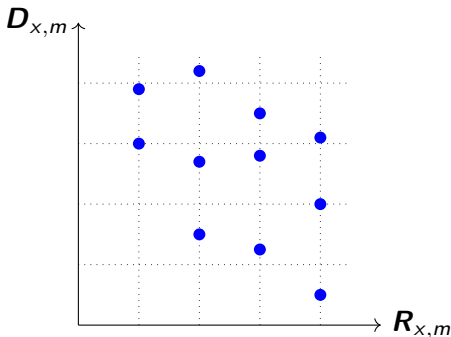
# Distortion-rate function: geometric solution via dual

- Dual:
$$D^*(R) = \sup_{\lambda \geq 0} \left\{ -\lambda R + \sum_{x \in \mathcal{X}} \min_{m \in \mathcal{M}_x} \left[ \boldsymbol{D}_{x,m} + \lambda \boldsymbol{R}_{x,m} \right] \right\}$$
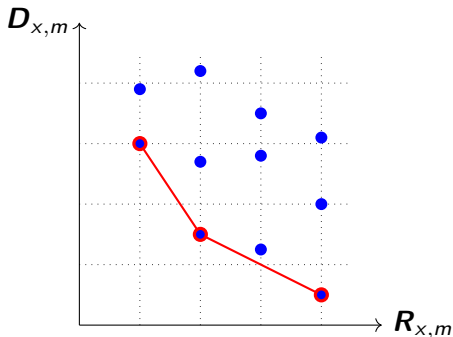
# Distortion-rate function: geometric solution via dual

- Dual:

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

# Distortion-rate function: geometric solution via dual

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

- Fix $\lambda \geq 0$, $x \in \mathcal{X}$
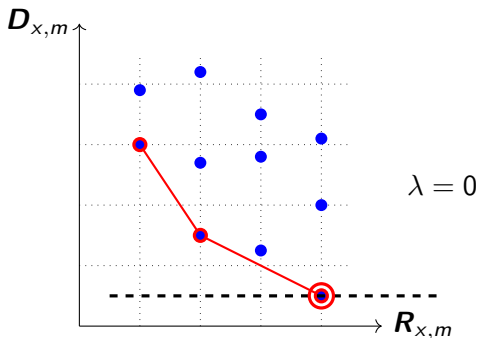
# Distortion-rate function: geometric solution via dual

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

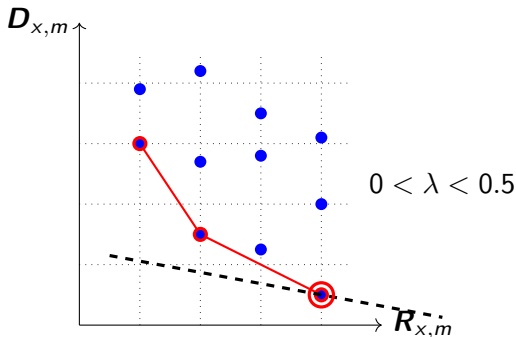- Fix $\lambda \geq 0$, $x \in \mathcal{X}$



- Relevant points: 32,000[10]

# Distortion-rate function: geometric solution via dual

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

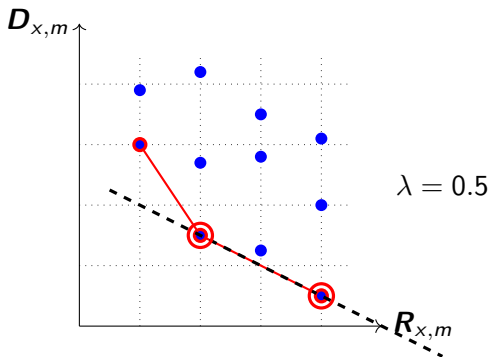- Fix $\lambda \geq 0$, $x \in \mathcal{X}$



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

# Distortion-rate function: geometric solution via dual

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

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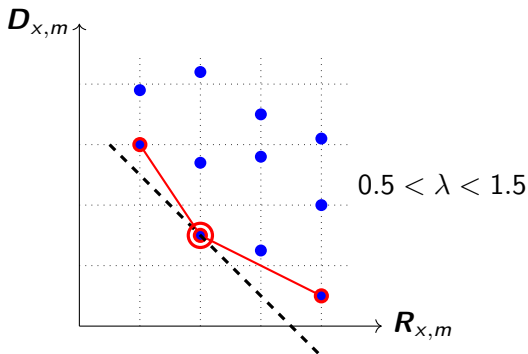


$\lambda = 0$

- Relevant points: 32,000[10]

# Distortion-rate function: geometric solution via dual

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

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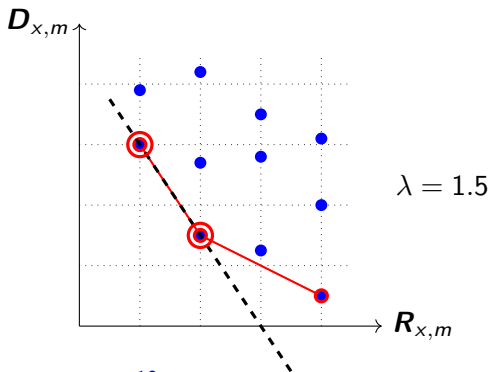


$0 < \lambda < 0.5$

- Relevant points: 32,000[10]

# Distortion-rate function: geometric solution via dual

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

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



$\lambda = 0.5$

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

# Distortion-rate function: geometric solution via dual

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

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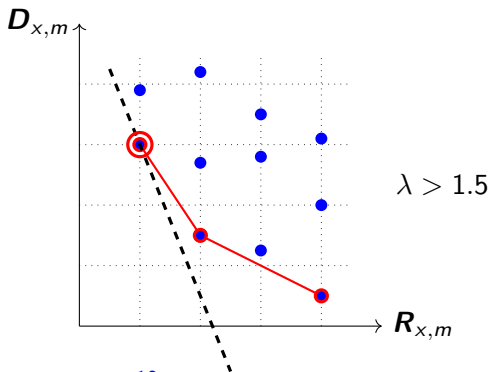


$0.5 < \lambda < 1.5$

- Relevant points: 32,000[10]

# Distortion-rate function: geometric solution via dual

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

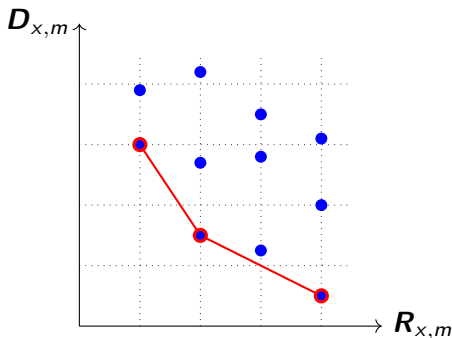- Fix $\lambda \geq 0$, $x \in \mathcal{X}$



$\lambda = 1.5$

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

# Distortion-rate function: geometric solution via dual

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

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



$\lambda > 1.5$
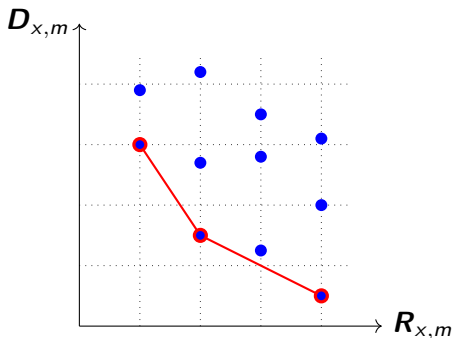
- Relevant points: $32{,}000^{10}$

# Distortion-rate function: geometric solution via dual

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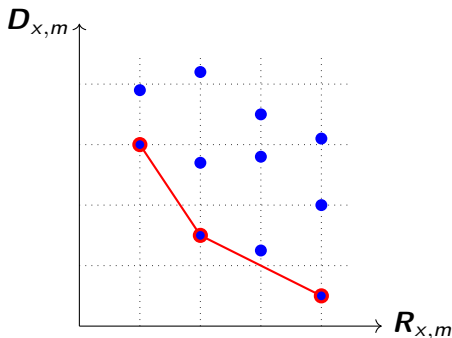
- Fix $\lambda \geq 0$, $x \in \mathcal{X}$



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

# Distortion-rate function: geometric solution via dual

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

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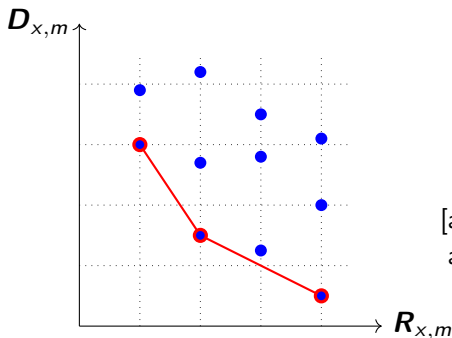


- Relevant points: $32{,}000^{10} \to 10$,
  only finitely many $\lambda$

# Distortion-rate function: geometric solution via dual

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

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



- Relevant points: $32{,}000^{10} \to 10$, $\qquad\qquad$ $(2^{10} \to 10)$
  only finitely many $\lambda$

# Distortion-rate function: geometric solution via dual

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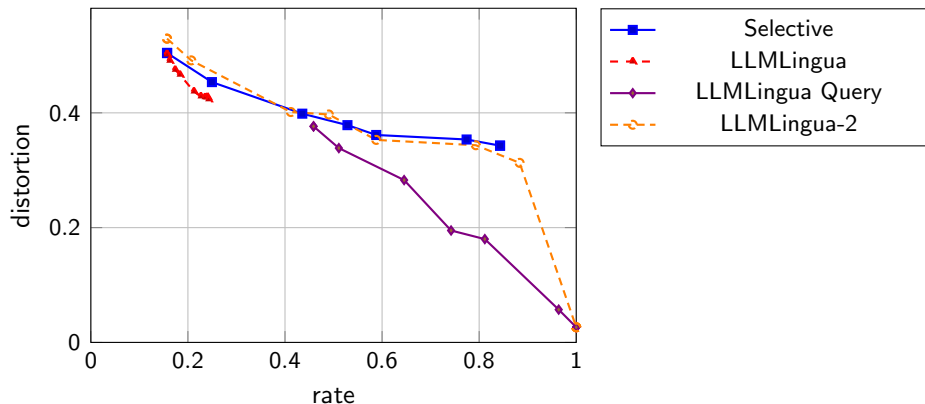
- Fix $\lambda \geq 0$, $x \in \mathcal{X}$



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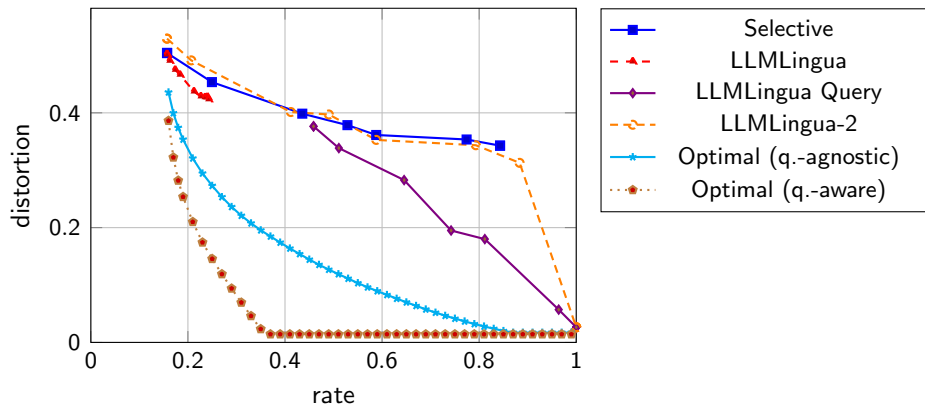apple $\not\mapsto$ pale, red, lp ]

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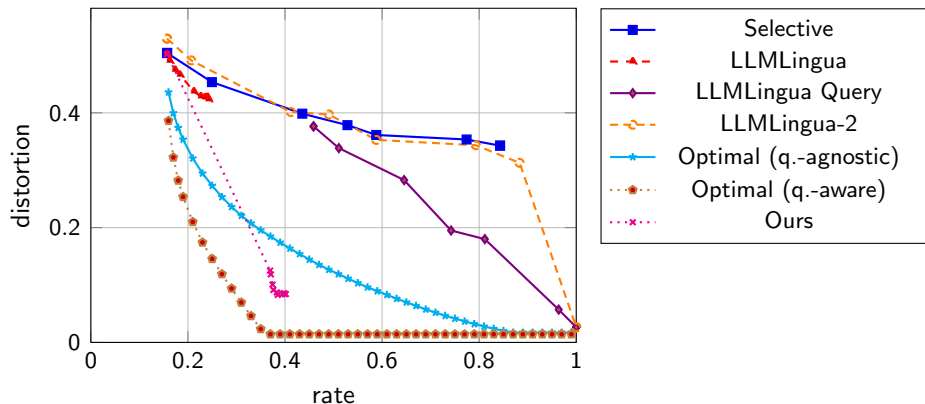  only finitely many $\lambda$

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

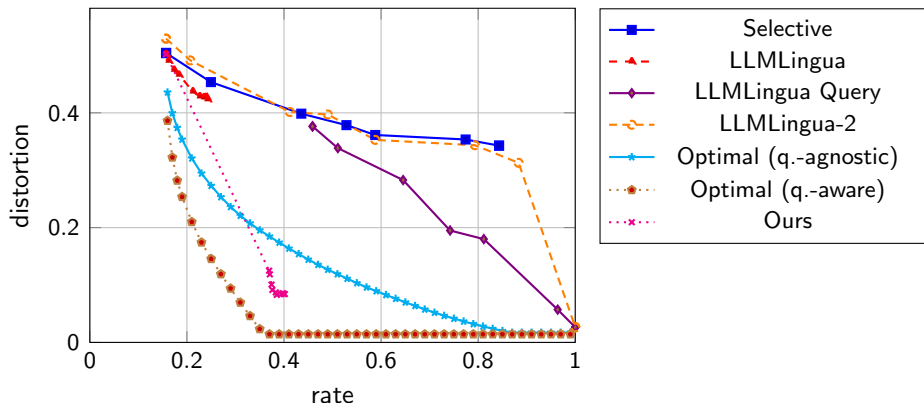# Experimental results

# Experimental results

# Experimental results
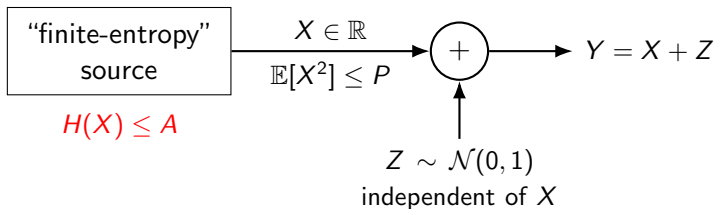
# Experimental results



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 [conference]

# Segue to a contraction problem

- Optimization 101...

# Segue to a contraction problem

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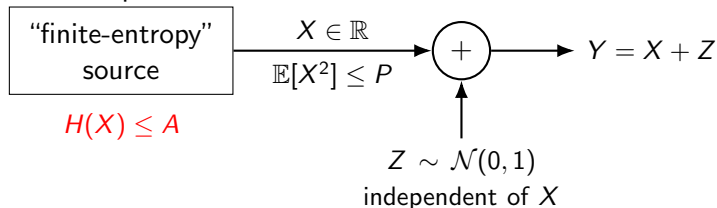


$$H(X) \leq A$$

# Outline

# Input-entropy-constrained channel capacity

- A contraction problem in communication



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

# 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



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

- Cardinality bounds? Finite support?

# Input-entropy-constrained channel capacity
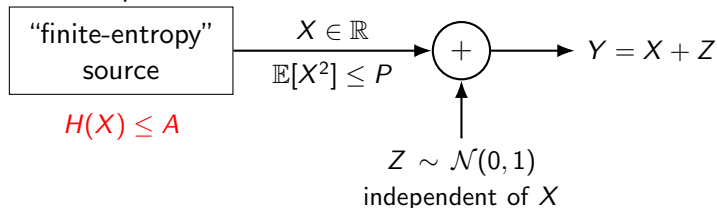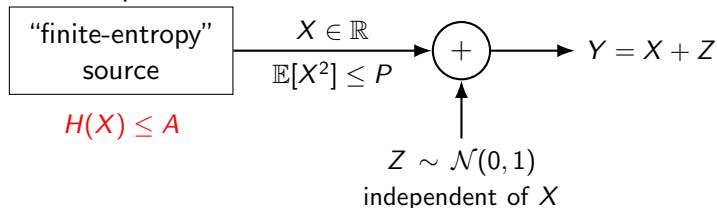
- A contraction problem in communication
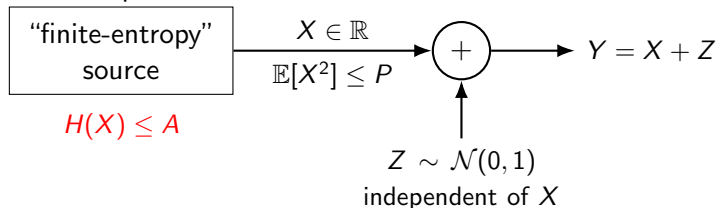


$$C_H(A, P) = \sup_{\substack{P_X: \; \mathbb{E}[X^2] \leq P \\ H(X) \leq A}} I(X; Y)$$

- Cardinality bounds? Finite support?

- A nontrivial upper bound better than

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

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*: $F_I(t) = \sup\limits_{P_{WX}: I(W;X) \leq t} I(W; Y)$

# Aside on data processing inequalities

Fix $P_{Y|X}$

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

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

# Aside on data processing inequalities

Fix $P_{Y|X}$

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

- Data processing *function*: $F_I(t) = \sup\limits_{P_{WX} : I(W;X) \leq t} I(W; Y)$

# Aside on data processing inequalities

Fix $P_{Y|X}$

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

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

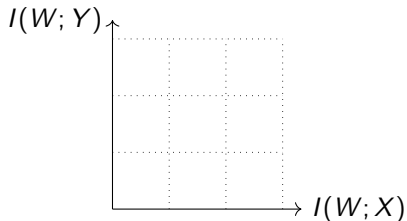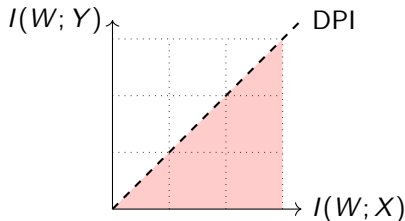# Aside on data processing inequalities
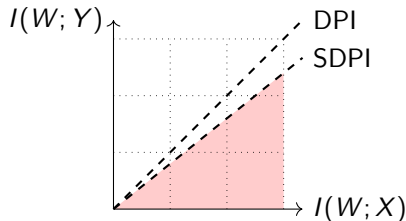
Fix $P_{Y|X}$

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

- Data processing *function*: $F_I(t) = \sup\limits_{P_{WX}: I(W;X) \leq t} I(W; Y)$



- 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*: $F_I(t) = \sup\limits_{P_{WX}:\, I(W;X) \leq t} I(W; Y)$



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

$$P_Y = P_X \circ P_{Y|X}$$

# Aside on data processing inequalities

Fix $P_{Y|X}$

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

- Data processing *function*: $F_I(t) = \sup\limits_{P_{WX} : I(W;X) \leq t} I(W; Y)$



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

- Natural analogue: $F_f(t) = \sup\limits_{P_X, Q_X : D_f(Q_X \| P_X) \leq t} D_f(Q_Y \| P_Y)$
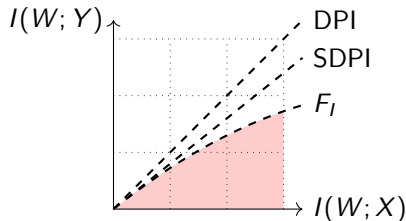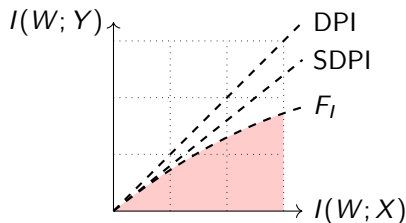
# Aside on data processing inequalities

Fix $P_{Y|X}$

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

- Data processing *function*: $F_I(t) = \sup\limits_{P_{WX}: I(W;X) \leq t} I(W; Y)$
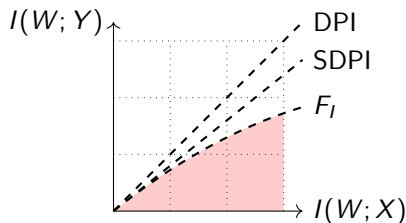


- Also DPI: for any $P_X, Q_X$, $D_f(Q_Y \,\|\, P_Y) \leq D_f(Q_X \,\|\, P_X)$

- Natural analogue: $F_f(t) = \sup\limits_{P_X, Q_X: D_f(Q_X \,\|\, P_X) \leq t} D_f(Q_Y \,\|\, P_Y)$

# Outline

Fix $P_{Y|X}$ and $f$

# Joint range of input and output divergences

Fix $P_{Y|X}$ and $f$

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

# Joint range of input and output divergences

Fix $P_{Y|X}$ and $f$

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

- Upper boundary of $\mathcal{D}_f = \bigcup_{P_X, Q_X} \left\{ \left( D_f(Q_X \| P_X), D_f(Q_Y \| P_Y) \right) \right\}$

# Joint range of input and output divergences

Fix $P_{Y|X}$ and $f$

- Define $F_f(t) = \sup\limits_{P_X, Q_X : D_f(Q_X \| P_X) \le t} D_f(Q_Y \| P_Y)$

- Upper boundary of $\mathcal{D}_f = \bigcup_{P_X, Q_X} \left\{ \left( D_f(Q_X \| P_X), D_f(Q_Y \| P_Y) \right) \right\}$

- Conjecture: $\mathcal{D}_f$ is convex                                    ?

# Joint range of input and output divergences

Fix $P_{Y|X}$ and $f$

- Define $F_f(t) = \sup\limits_{P_X, Q_X \colon D_f(Q_X \,\|\, P_X) \leq t} D_f(Q_Y \,\|\, P_Y)$

- Upper boundary of $\mathcal{D}_f = \bigcup_{P_X, Q_X} \left\{ \left( D_f(Q_X \,\|\, P_X), D_f(Q_Y \,\|\, P_Y) \right) \right\}$

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

- Upper boundary of $\mathcal{D}_f = \bigcup_{P_X, Q_X} \left\{ \left( D_f(Q_X \,\|\, P_X), D_f(Q_Y \,\|\, P_Y) \right) \right\}$

- 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 $F_f(t) = \sup\limits_{P_X, Q_X : D_f(Q_X \| P_X) \le t} D_f(Q_Y \| P_Y)$

- Upper boundary of $\mathcal{D}_f = \bigcup_{P_X, Q_X} \left\{ \left( D_f(Q_X \| P_X), D_f(Q_Y \| P_Y) \right) \right\}$

- Conjecture: $\mathcal{D}_f$ is convex ( $\implies F_f$ is concave)      ?

- Facts:
  - $F_I$ is NOT necessarily concave (counter-example: $P_{Y|X} = BEC^3$)

# Joint range of input and output divergences

Fix $P_{Y|X}$ and $f$

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

- Upper boundary of $\mathcal{D}_f = \bigcup\limits_{P_X, Q_X} \left\{ \left( D_f(Q_X \| P_X), D_f(Q_Y \| P_Y) \right) \right\}$

- Conjecture: $\mathcal{D}_f$ is convex ( $\implies F_f$ is concave) ?

- Facts:
    - $F_I$ is NOT necessarily concave (counter-example: $P_{Y|X} = \mathrm{BEC}^3$)
    - Fix $P_X$, define $\tilde{F}_I(t, P_X) = \sup\limits_{P_{W|X} : I(W;X) \leq t} I(W;Y)$ and

# Joint range of input and output divergences

Fix $P_{Y|X}$ and $f$

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

- Upper boundary of $\mathcal{D}_f = \bigcup_{P_X, Q_X} \left\{ \left( D_f(Q_X \| P_X), D_f(Q_Y \| P_Y) \right) \right\}$

- 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 $\tilde{F}_I(t, P_X) = \sup\limits_{P_{W|X} : I(W;X) \leq t} I(W;Y)$ and
    $$\tilde{F}_f(t, P_X) = \sup\limits_{Q_X : D_f(Q_X \| P_X) \leq t} D_f(Q_Y \| P_Y),$$

# Joint range of input and output divergences

Fix $P_{Y|X}$ and $f$

- Define $F_f(t) = \sup\limits_{P_X, Q_X \colon D_f(Q_X \| P_X) \leq t} D_f(Q_Y \| P_Y)$

- Upper boundary of $\mathcal{D}_f = \bigcup_{P_X, Q_X} \left\{ \left( D_f(Q_X \| P_X), D_f(Q_Y \| P_Y) \right) \right\}$

- Conjecture: $\mathcal{D}_f$ is convex ( $\implies$ $F_f$ is concave)       ?

- Facts:
    - $F_I$ is NOT necessarily concave (counter-example: $P_{Y|X} = \mathrm{BEC}^3$)

    - Fix $P_X$, define $\tilde{F}_I(t, P_X) = \sup\limits_{P_{W|X} \colon I(W;X) \leq t} I(W;Y)$ and
    $$\tilde{F}_f(t, P_X) = \sup\limits_{Q_X \colon D_f(Q_X \| P_X) \leq t} D_f(Q_Y \| P_Y),$$
    then $\tilde{F}_I(\cdot, P_X)$ is concave

# Joint range of input and output divergences

Fix $P_{Y|X}$ and $f$

- Define $F_f(t) = \sup\limits_{P_X, Q_X \colon D_f(Q_X \| P_X) \leq t} D_f(Q_Y \| P_Y)$

- Upper boundary of $\mathcal{D}_f = \bigcup\limits_{P_X, Q_X} \left\{ \left( D_f(Q_X \| P_X), D_f(Q_Y \| P_Y) \right) \right\}$

- 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 $\tilde{F}_I(t, P_X) = \sup\limits_{P_{W|X} \colon I(W;X) \leq t} I(W;Y)$ and
      $$\tilde{F}_f(t, P_X) = \sup\limits_{Q_X \colon D_f(Q_X \| P_X) \leq t} D_f(Q_Y \| P_Y),$$
      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 $F_f(t) = \sup\limits_{P_X, Q_X : D_f(Q_X \| P_X) \le t} D_f(Q_Y \| P_Y)$

- Upper boundary of $\mathcal{D}_f = \bigcup_{P_X, Q_X} \left\{ \left( D_f(Q_X \| P_X), D_f(Q_Y \| P_Y) \right) \right\}$

- 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 $\tilde{F}_I(t, P_X) = \sup\limits_{P_{W|X} : I(W;X) \le t} I(W;Y)$ and
    $$\tilde{F}_f(t, P_X) = \sup\limits_{Q_X : D_f(Q_X \| P_X) \le t} D_f(Q_Y \| P_Y),$$
    then $\tilde{F}_I(\cdot, P_X)$ is concave; $\implies \tilde{F}_f(\cdot, P_X)$ is concave

    - For any $f, g$, $\bigcup_{P_X, Q_X} \left\{ \left( D_f(Q_X \| P_X), D_g(Q_X \| P_X) \right) \right\}$ is convex

# Outline

# 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 $\longrightarrow$ 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 $\longrightarrow$ 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!