



# Sketch In, Sketch Out: Accelerating both Learning and Inference for Structured Prediction with Kernels

AISTATS 2024

---

Tamim El Ahmad<sup>\*</sup>, Luc Brogat-Motte<sup>\*†</sup>, Pierre Laforgue<sup>‡</sup>, Florence d'Alché-Buc<sup>\*</sup>

<sup>\*</sup> LTCI, Télécom Paris, Institut Polytechnique de Paris

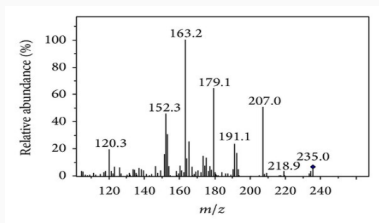
<sup>†</sup> L2S, CentraleSupélec

<sup>‡</sup> Università degli Studi di Milano

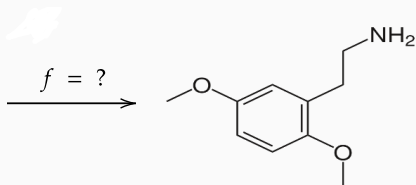
July 16, 2024

# Structured prediction

Emblematic example of metabolite identification (Brouard et al., 2016a; Schymanski et al., 2017):



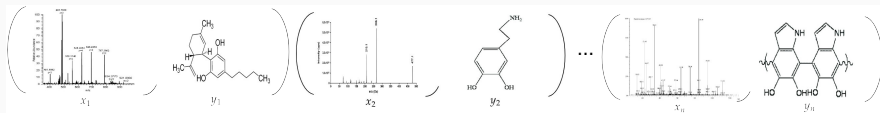
$x$



$y$

# Structured prediction in supervised settings

Supervised settings:  $n$  i.i.d. training sample  $(x_i, y_i)_{i=1}^n \in (\mathcal{X}, \mathcal{Y})^n \sim \rho$

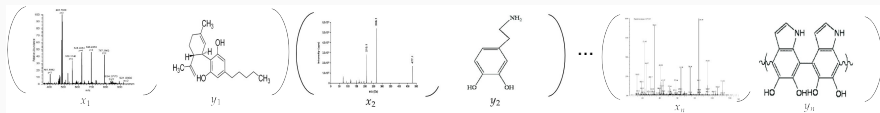


Given a loss function  $\Delta : \mathcal{Y}^2 \rightarrow \mathbb{R}$

$$f^* = \arg \inf_{f: \mathcal{X} \rightarrow \mathcal{Y}} \mathbb{E}_{(x,y) \sim \rho} [\Delta(f(x), y)] \approx \arg \inf_{f: \mathcal{X} \rightarrow \mathcal{Y}} \frac{1}{n} \sum_{i=1}^n \Delta(f(x_i), y_i) = \hat{f}$$

# Structured prediction in supervised settings

Supervised settings:  $n$  i.i.d. training sample  $(x_i, y_i)_{i=1}^n \in (\mathcal{X}, \mathcal{Y})^n \sim \rho$



Given a loss function  $\Delta : \mathcal{Y}^2 \rightarrow \mathbb{R}$

$$f^* = \arg \inf_{f: \mathcal{X} \rightarrow \mathcal{Y}} \mathbb{E}_{(x,y) \sim \rho} [\Delta(f(x), y)] \approx \arg \inf_{f: \mathcal{X} \rightarrow \mathcal{Y}} \frac{1}{n} \sum_{i=1}^n \Delta(f(x_i), y_i) = \hat{f}$$

How to design a loss  $\Delta$  taking into account the structure of  $\mathcal{Y}$ ?

# Table of contents

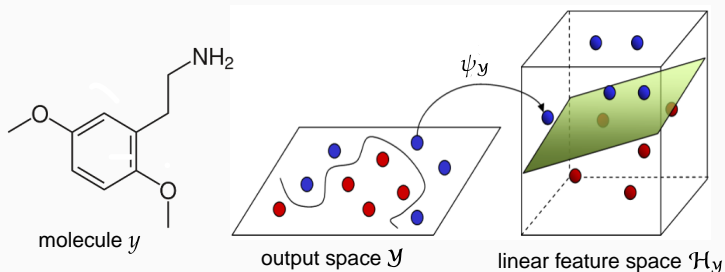
1. Input Output Kernel Regression
2. Sketched Input Sketched Output Kernel Regression
3. Theoretical guarantees
4. Experiments
5. Conclusion

# Input Output Kernel Regression

---

# Kernel methods: output representation

Linear method after embedding through feature map  $\psi_{\mathcal{Y}} : \mathcal{Y} \rightarrow \mathcal{H}_{\mathcal{Y}}$ :  
choice of kernel  $\iff$  choice of representation



$\langle \psi_{\mathcal{Y}}(y), \psi_{\mathcal{Y}}(y') \rangle_{\mathcal{H}_{\mathcal{Y}}} = k_{\mathcal{Y}}(y, y')$ : relevant similarity measure over  $\mathcal{Y}$

$$\implies \Delta(y, y') = \|\psi_{\mathcal{Y}}(y) - \psi_{\mathcal{Y}}(y')\|_{\mathcal{H}_{\mathcal{Y}}}^2 = 2 - 2k_{\mathcal{Y}}(y, y')$$

( $\forall y \in \mathcal{Y}, \|\psi_{\mathcal{Y}}\|_{\mathcal{H}_{\mathcal{Y}}} = 1$  without loss of generality)

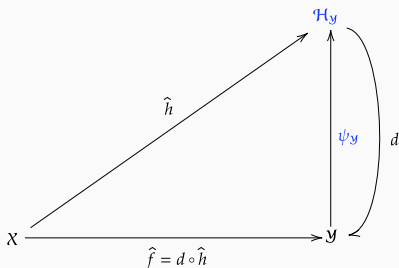
**Versatility:** tackle various tasks through an appropriate choice of  $\psi_{\mathcal{Y}}$  (e.g. SOTA performance on metabolite identification (Brouard et al., 2016a) and label ranking (Korba et al., 2018) datasets)



# Output Kernel Regression: a surrogate approach

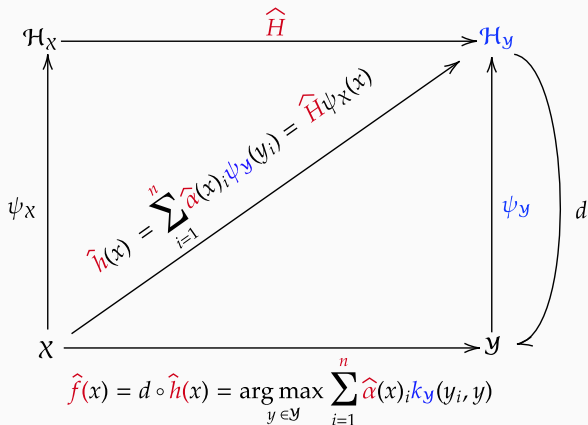
Surrogate (2-step) method (Weston et al., 2003; Cortes et al., 2005; Brouard et al., 2011; Kadri et al., 2013):

1.  $\hat{h} = \arg \min_{h: \mathcal{X} \rightarrow \mathcal{H}_Y} \frac{1}{n} \sum_{i=1}^n \|h(x_i) - \psi_Y(y_i)\|_{\mathcal{H}_Y}^2$  (training step)
2.  $\hat{f}(x) = d \circ \hat{h}(x) = \arg \min_{y \in \mathcal{Y}} \|\hat{h}(x) - \psi_Y(y)\|_{\mathcal{H}_Y}^2$  (inference step)



**Theoretical guarantees:** for measurable  $h : \mathcal{X} \rightarrow \mathcal{H}_Y$  and  $f = d \circ h$ ,  $\hat{f}$ 's excess risk is bounded by  $\hat{h}$ 's excess risk (Ciliberto et al., 2020)

# Input Output Kernel Regression



**IOKR:** Weston et al. (2003); Cortes et al. (2005); Brouard et al. (2011); Kadri et al. (2013); Brouard et al. (2016b); Korba et al. (2018)

# IOKR: training and inference complexities

1. Training:  $\hat{h}(x) = \sum_{i=1}^n \hat{\alpha}(x)_i \psi_{\mathcal{Y}}(y_i)$  where

$$\hat{\alpha}(x) = \underbrace{(K_X + n\lambda I_n)}_{n \times n}^{-1} k_X^x = \hat{\Omega} k_X^x$$

$\Rightarrow \mathcal{O}(n^3)$  time complexity

# IOKR: training and inference complexities

1. Training:  $\hat{h}(x) = \sum_{i=1}^n \hat{\alpha}(x)_i \psi_{\mathcal{Y}}(y_i)$  where

$$\hat{\alpha}(x) = \underbrace{(K_X + n\lambda I_n)}_{n \times n}^{-1} k_X^x = \hat{\Omega} k_X^x$$

$\Rightarrow \mathcal{O}(n^3)$  time complexity

2. Inference:  $\hat{f}(x) = \arg \max_{y \in \mathcal{Y}} \sum_{i=1}^n \hat{\alpha}(x)_i k_{\mathcal{Y}}(y_i, y) = k_X^{xT} \hat{\Omega} k_Y^y$

- Test set:  $X^{\text{te}} = \{x_1^{\text{te}}, \dots, x_{n_{\text{te}}}^{\text{te}}\}$  of size  $n_{\text{te}}$
- Candidate set:  $Y^{\text{c}} = \{y_1^{\text{c}}, \dots, y_{n_{\text{c}}}^{\text{c}}\}$  of size  $n_{\text{c}}$

$$\underbrace{K_X^{\text{te, tr}}}_{n_{\text{te}} \times n} \underbrace{\hat{\Omega}}_{n \times n} \underbrace{K_Y^{\text{tr, c}}}_{n \times n_{\text{c}}}$$

$$\hat{f}(x_i^{\text{te}}) = y_j^{\text{c}} \quad \text{where} \quad j = \arg \max_{1 \leq j \leq n_{\text{c}}} [K_X^{\text{te, tr}} \hat{\Omega} K_Y^{\text{tr, c}}]_{ij}$$

$\Rightarrow \mathcal{O}(n_{\text{te}} n n_{\text{c}})$  time complexity if  $n_{\text{te}} < n \leq n_{\text{c}}$

1. **Scalability:** obtain  $\tilde{f} = d \circ \tilde{h}$ , **computationally efficient** version of  $\hat{f} = d \circ \hat{h}$ , when learning from **big data**, i.e. **large  $n$**

2. **Theory:** obtain **excess risk bound** of  $\tilde{f} = d \circ \tilde{h}$

# Key tool for scalability: Random Fourier Features vs Sketching

a) Random Fourier Features (Rahimi and Recht, 2007; Sriperumbudur and Szabó, 2015): for  $m_{\mathcal{Y}} \ll n$ ,

$$\langle \psi_{\mathcal{Y}}(y), \psi_{\mathcal{Y}}(y') \rangle_{\mathcal{H}_{\mathcal{Y}}} \approx \langle \tilde{\psi}_{\mathcal{Y}}(y), \tilde{\psi}_{\mathcal{Y}}(y') \rangle_{\mathbb{R}^{m_{\mathcal{Y}}}}$$

$$\implies \Delta(y, y') = \|\psi_{\mathcal{Y}}(y) - \psi_{\mathcal{Y}}(y')\|_{\mathcal{H}_{\mathcal{Y}}}^2 \approx \|\tilde{\psi}_{\mathcal{Y}}(y) - \tilde{\psi}_{\mathcal{Y}}(y')\|_{\mathbb{R}^{m_{\mathcal{Y}}}}^2 = \tilde{\Delta}(y, y')$$

$$\implies \tilde{\Delta} \text{ approximated loss}$$

# Key tool for scalability: Random Fourier Features vs Sketching

a) **Random Fourier Features** (Rahimi and Recht, 2007; Sriperumbudur and Szabó, 2015): for  $m_{\mathcal{Y}} \ll n$ ,

$$\langle \psi_{\mathcal{Y}}(y), \psi_{\mathcal{Y}}(y') \rangle_{\mathcal{H}_{\mathcal{Y}}} \approx \langle \tilde{\psi}_{\mathcal{Y}}(y), \tilde{\psi}_{\mathcal{Y}}(y') \rangle_{\mathbb{R}^{m_{\mathcal{Y}}}}$$

$$\implies \Delta(y, y') = \|\psi_{\mathcal{Y}}(y) - \psi_{\mathcal{Y}}(y')\|_{\mathcal{H}_{\mathcal{Y}}}^2 \approx \|\tilde{\psi}_{\mathcal{Y}}(y) - \tilde{\psi}_{\mathcal{Y}}(y')\|_{\mathbb{R}^{m_{\mathcal{Y}}}}^2 = \tilde{\Delta}(y, y')$$

$\implies \tilde{\Delta}$  approximated loss

b) **Sketching** (Williams and Seeger, 2001; Rudi et al., 2015; Yang et al., 2017): for  $m_{\mathcal{Y}} \ll n$ ,  $R_{\mathcal{Y}} \in \mathbb{R}^{m_{\mathcal{Y}} \times n}$

$$\text{span} \left( (\psi_{\mathcal{Y}}(y_i))_{i=1}^n \right) \leftarrow \text{span} \left( \left( \sum_{j=1}^n [R_{\mathcal{Y}}]_{ij} \psi_{\mathcal{Y}}(y_j) \right)_{i=1}^{m_{\mathcal{Y}}} \right)$$

$\implies \Delta$  remains unchanged!

# Sketched Input Sketched Output Kernel Regression

---



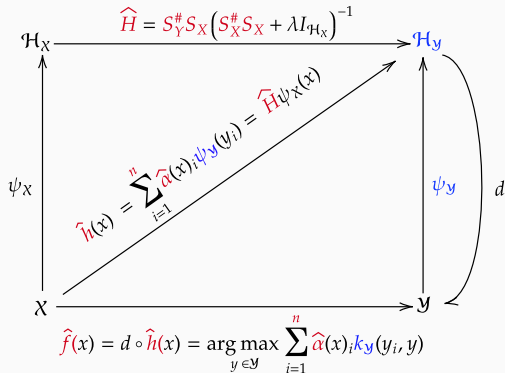
**Motivation:** build a **low-rank** approximation  $\tilde{h}$  of  $\hat{h}$  thanks to **input and output** random projectors  $\tilde{P}_X$  and  $\tilde{P}_Y$  to obtain a **scalable** predictor  $\tilde{f}$  together with an **excess risk bound**

## Some notations

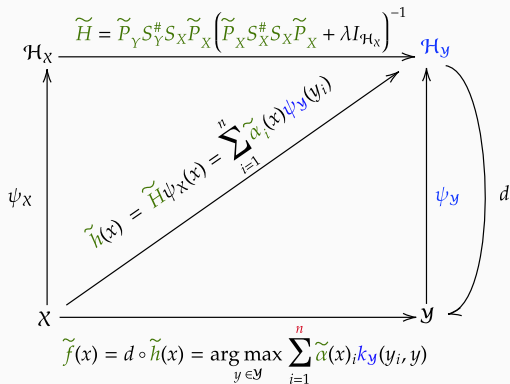
For an i.i.d. sample  $(z_i)_{i=1}^n \in \mathcal{Z}^n \sim \rho_{\mathcal{Z}}$ :

- $S_{\mathcal{Z}} : f \in \mathcal{H}_{\mathcal{Z}} \mapsto (1/\sqrt{n})(\langle f, \psi_{\mathcal{Z}}(z_1) \rangle_{\mathcal{H}_{\mathcal{Z}}}, \dots, \langle f, \psi_{\mathcal{Z}}(z_n) \rangle_{\mathcal{H}_{\mathcal{Z}}})^{\top} \in \mathbb{R}^n$   
sampling operator
- $S_{\mathcal{Z}}^{\#} : \alpha \in \mathbb{R}^n \mapsto (1/\sqrt{n}) \sum_{i=1}^n \alpha_i \psi_{\mathcal{Z}}(z_i) \in \text{span}((\psi_{\mathcal{Z}}(z_i))_{i=1}^n)$  its  
adjoint
- $C_{\mathcal{Z}} = \mathbb{E}_{\mathcal{Z}}[\psi_{\mathcal{Z}}(z) \otimes \psi_{\mathcal{Z}}(z)]$  covariance operator
- $\widehat{C}_{\mathcal{Z}} = (1/n) \sum_{i=1}^n \psi_{\mathcal{Z}}(z_i) \otimes \psi_{\mathcal{Z}}(z_i) = S_{\mathcal{Z}}^{\#} S_{\mathcal{Z}}$  its empirical counterpart:  
 $\widehat{C}_{\mathcal{Z}} : \mathcal{H}_{\mathcal{Z}} \rightarrow \text{span}((\psi_{\mathcal{Z}}(z_i))_{i=1}^n)$

# Low-rank estimator: from IOKR to SISOKR



# Low-rank estimator: from IOKR to SISOKR



$$\tilde{P}_Z : \mathcal{H}_Z \rightarrow \tilde{\mathcal{H}}_Z \text{ where } \tilde{\mathcal{H}}_Z := \text{span} \left( \left( \sum_{j=1}^n [\mathbf{R}_Z]_{ij} \psi_Z(z_j) \right)_{i=1}^{m_Z} \right)$$

How to build these projectors?

## Construction of the orthogonal projector $\tilde{P}_Z$

- $\hat{C}_Z = S_Z^\# S_Z = (1/n) \sum_{i=1}^n \psi_Z(z_i) \otimes \psi_Z(z_i)$
- $\tilde{C}_Z = S_Z^\# R_Z^\top R_Z S_Z = \frac{1}{n} \sum_{l=1}^{m_Z} \left( \sum_{i=1}^n R_{z_{li}} \psi_Z(z_i) \right) \otimes \left( \sum_{j=1}^n R_{z_{lj}} \psi_Z(z_j) \right)$
- $\tilde{K}_Z = R_Z K_Z R_Z^\top$ , and  $\left\{ \left( \sigma_i(\tilde{K}_Z), \tilde{u}_i^Z \right), i \in [m_Z] \right\}$  its eigenpairs
- $p_Z = \text{rank}(\tilde{K}_Z)$ , and for all  $1 \leq i \leq p_Z$ ,  $\tilde{e}_i^Z = \sqrt{\frac{n}{\sigma_i(\tilde{K}_Z)}} S_Z^\# R_Z^\top \tilde{u}_i^Z \in \mathcal{H}_Z$

# Construction of the orthogonal projector $\tilde{P}_Z$

- $\hat{C}_Z = S_Z^\# S_Z = (1/n) \sum_{i=1}^n \psi_Z(z_i) \otimes \psi_Z(z_i)$
- $\tilde{C}_Z = S_Z^\# R_Z^\top R_Z S_Z = \frac{1}{n} \sum_{l=1}^{m_Z} \left( \sum_{i=1}^n R_{Z_{li}} \psi_Z(z_i) \right) \otimes \left( \sum_{j=1}^n R_{Z_{lj}} \psi_Z(z_j) \right)$
- $\tilde{K}_Z = R_Z K_Z R_Z^\top$ , and  $\left\{ \left( \sigma_i(\tilde{K}_Z), \tilde{u}_i^Z \right), i \in [m_Z] \right\}$  its eigenpairs
- $p_Z = \text{rank}(\tilde{K}_Z)$ , and for all  $1 \leq i \leq p_Z$ ,  $\tilde{e}_i^Z = \sqrt{\frac{n}{\sigma_i(\tilde{K}_Z)}} S_Z^\# R_Z^\top \tilde{u}_i^Z \in \mathcal{H}_Z$

Proposition (El Ahmad et al., 2024)

The  $\tilde{e}_i^Z$ s are the **eigenfunctions**, associated to the eigenvalues  $\sigma_i(\tilde{K}_Z)/n$ , of  $\tilde{C}_Z$ , whose range is  $\text{span}((\sum_{j=1}^n R_{Z_{ij}} \psi_Z(z_j))_{i=1}^{m_Z})$ . Then,  $\tilde{E}^Z = (\tilde{e}_1^Z, \dots, \tilde{e}_{p_Z}^Z)$  is an **orthonormal basis** of  $\text{span}((\sum_{j=1}^n R_{Z_{ij}} \psi_Z(z_j))_{i=1}^{m_Z})$ , and  $\tilde{P}_Z$  writes as

$$\tilde{P}_Z = \sum_{i=1}^{p_Z} \langle \cdot, \tilde{e}_i^Z \rangle_{\mathcal{H}_Z} \tilde{e}_i^Z = (R_Z S_Z)^\# (R_Z S_Z (R_Z S_Z)^\#)^\dagger R_Z S_Z.$$

Proposition (El Ahmad et al., 2024)

$$\tilde{h}(x) = \sum_{i=1}^n \tilde{\alpha}_i(x) \psi_Y(y_i), \quad \text{where} \quad \tilde{\alpha}(x) = R_Y^\top \tilde{\Omega} R_X k_X^x,$$

with

$$\tilde{\Omega} = \underbrace{(R_Y K_Y R_Y^\top)^\dagger}_{m_Y \times m_Y} R_Y K_Y K_X R_X^\top \underbrace{(R_X K_X^2 R_X^\top + n\lambda R_X K_X R_X^\top)^\dagger}_{m_X \times m_X}$$

Proposition (El Ahmad et al., 2024)

$$\tilde{h}(x) = \sum_{i=1}^n \tilde{\alpha}_i(x) \psi_y(y_i), \quad \text{where} \quad \tilde{\alpha}(x) = R_y^\top \tilde{\Omega} R_x k_x^x,$$

with

$$\tilde{\Omega} = \underbrace{(R_y K_y R_y^\top)^\dagger}_{m_y \times m_y} R_y K_y K_x R_x^\top \underbrace{(R_x K_x^2 R_x^\top + n \lambda R_x K_x R_x^\top)^\dagger}_{m_x \times m_x}$$

Inversion complexity:  $\mathcal{O}(n^3) \rightarrow \mathcal{O}(\max(m_x^3, m_y^3))$

Complexity of  $R_Z K_Z$ : depends on the sketching matrix, between  $\mathcal{O}(nm_Z)$  and  $\mathcal{O}(n^2 m_Z)$

$\implies$  Training complexity reduced thanks to input sketching!



$$\tilde{f}(x) = \arg \max_{y \in \mathcal{Y}} \sum_{i=1}^n \tilde{\alpha}(x)_i k_{\mathcal{Y}}(y_i, y) = \arg \max_{y \in \mathcal{Y}} k_X^{x,T} R_X^T \tilde{\Omega} R_Y k_Y^y$$

$$\underbrace{K_X^{\text{te, tr}} R_X^T}_{n_{\text{te}} \times m_X} \underbrace{\tilde{\Omega}}_{m_X \times m_Y} \underbrace{R_Y K_Y^{\text{tr, c}}}_{m_Y \times n_c}$$

$$\tilde{f}(x_i^{\text{te}}) = y_j^c \quad \text{where} \quad j = \arg \max_{1 \leq j \leq n_c} [K_X^{\text{te, tr}} R_X^T \tilde{\Omega} R_Y K_Y^{\text{tr, c}}]_{ij}$$

$$\tilde{f}(x) = \arg \max_{y \in \mathcal{Y}} \sum_{i=1}^n \tilde{\alpha}(x)_i k_{\mathcal{Y}}(y_i, y) = \arg \max_{y \in \mathcal{Y}} k_X^{x^T} R_X^T \tilde{\Omega} R_Y k_Y^y$$

$$\underbrace{K_X^{\text{te, tr}} R_X^T}_{n_{\text{te}} \times m_X} \underbrace{\tilde{\Omega}}_{m_X \times m_Y} \underbrace{R_Y K_Y^{\text{tr, c}}}_{m_Y \times n_C}$$

$$\tilde{f}(x_i^{\text{te}}) = y_j^c \quad \text{where} \quad j = \arg \max_{1 \leq j \leq n_C} [K_X^{\text{te, tr}} R_X^T \tilde{\Omega} R_Y K_Y^{\text{tr, c}}]_{ij}$$

Decoding complexity:  $\mathcal{O}(n_{\text{te}} n n_C) \rightarrow \mathcal{O}(n_{\text{te}} m_Y n_C)$  if

$$n_{\text{te}} \leq m_X, m_Y < n \leq n_C$$

$\Rightarrow$  Inference complexity reduced thanks to output sketching!

Scalability  $\checkmark!$

# Theoretical guarantees

---

# Theoretical guarantees of SISOKR

Let

$$\mathcal{R}(f) = \mathbb{E}_{(x,y) \sim \rho} [\Delta(f(x), y)],$$

and

$$f^* = \arg \inf_{f: \mathcal{X} \rightarrow \mathcal{Y}} \mathbb{E}_{(x,y) \sim \rho} [\Delta(f(x), y)],$$

we want to control

$$\mathcal{R}(\tilde{f}) - \mathcal{R}(f^*) \leq ?$$

# Assumptions

**Asm. 1 (Attainability):** Recall that  $h^*(x) := \mathbb{E}_Y[\psi_Y(Y) \mid X = x]$ . There exists  $H : \mathcal{H}_X \rightarrow \mathcal{H}_Y$  with  $\|H\|_{\text{HS}} < +\infty$  such that

$$h^*(x) = H\psi_X(x) \quad \forall x \in \mathcal{X}.$$

**Asm. 2 (Bounded kernel):** there exists  $\kappa_Z > 0$  such that

$$k_Z(z, z) \leq \kappa_Z^2 \quad \forall z \in \mathcal{Z}.$$

**Asm. 3 (Capacity condition):** there exists  $\gamma_Z \in [0, 1]$  such that

$$Q_Z := \text{Tr}(C_Z^{\gamma_Z}) < +\infty.$$

**Asm. 4 (Embedding property):** there exists  $b_Z > 0$  and  $\mu_Z \in [0, 1]$  such that almost surely

$$\psi_Z(z) \otimes \psi_Z(z) \preceq b_Z C_Z^{1-\mu_Z}.$$

**Asm. 5 (Sub-Gaussian sketches):**  $R_Z \in \mathbb{R}^{m_Z \times n}$  composed with i.i.d. entries s.t. (i)  $\mathbb{E}[R_{Z,ij}] = 0$ , (ii)  $\mathbb{E}[R_{Z,ij}^2] = 1/m_Z$  and (iii)

$R_{Z,ij} \sim \frac{\nu_Z}{m_Z} - \text{sub-Gaussian with } \nu_Z \geq 1.$

# Theorem: SISOKR learning rates (El Ahmad et al., 2024)

Under Asm. 1, 2, 3, 4 and 5, if for all  $y \in \mathcal{Y}$ ,  $\|\psi_y(y)\|_{\mathcal{H}_y} = \kappa_y$ , for  $\mathcal{Z} \in \{\mathcal{X}, \mathcal{Y}\}$  and for  $n \in \mathbb{N}$  sufficiently large such that  $\frac{9}{n} \log(n/\delta) \leq n^{-\frac{1}{1+\gamma_{\mathcal{Z}}}} \leq \|C_{\mathcal{Z}}\|_{\text{op}}/2$ , and for sketching sizes  $m_{\mathcal{Z}}, \in \mathbb{N}$  such that

$$m_{\mathcal{Z}} \gtrsim \max \left( \nu_{\mathcal{Z}}^2 n^{\frac{\gamma_{\mathcal{Z}} + \mu_{\mathcal{Z}}}{1 + \gamma_{\mathcal{Z}}}}, \nu_{\mathcal{Z}}^4 \log(1/\delta) \right),$$

then with probability  $1 - \delta$

$$\mathbb{E}[\|\tilde{h}(x) - h^*(x)\|_{\mathcal{H}_y}^2]^{\frac{1}{2}} \lesssim \log(4/\delta) n^{-\frac{1 - \gamma_{\mathcal{X}} \vee \gamma_{\mathcal{Y}}}{2(1 + \gamma_{\mathcal{X}} \vee \gamma_{\mathcal{Y}})}},$$

and

$$\mathcal{R}(\tilde{f}) - \mathcal{R}(f^*) \lesssim \mathbb{E}[\|\tilde{h}(x) - h^*(x)\|_{\mathcal{H}_y}^2]^{\frac{1}{2}} \lesssim \log(4/\delta) n^{-\frac{1 - \gamma_{\mathcal{X}} \vee \gamma_{\mathcal{Y}}}{2(1 + \gamma_{\mathcal{X}} \vee \gamma_{\mathcal{Y}})}}.$$

# Experiments

---

# Synthetic least squares regression

1)  $n = 10\,000$ ,  $\mathcal{X} = \mathcal{Y} = \mathbb{R}^d$ ,  $d = 300$ ,  $k_{\mathcal{X}}$  and  $k_{\mathcal{Y}}$  linear kernels  $\implies$   
 $\mathcal{H}_{\mathcal{X}} = \mathcal{H}_{\mathcal{Y}} = \mathbb{R}^d$

2) Construct covariance matrices  $C_{\mathcal{X}}$  and  $E$  such that  $\sigma_k(C_{\mathcal{X}}) = k^{-3/2}$   
and  $\sigma_k(E) = 0.2k^{-1/10}$

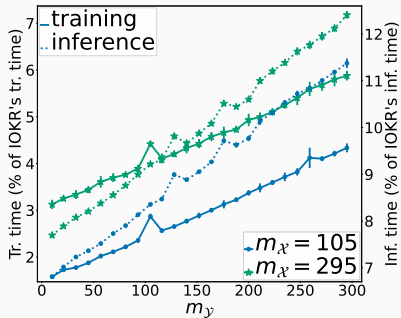
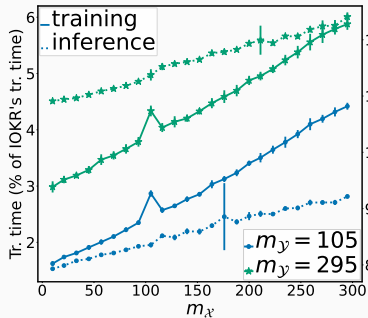
3) Draw  $H_0 \sim \mathcal{N}(0, I_d)$ , and for  $i \leq n$ ,  $x_i \sim \mathcal{N}(0, C_{\mathcal{X}})$ ,  $\epsilon_i \sim \mathcal{N}(0, E)$ ,

$$y_i = C_{\mathcal{X}} H_0 x_i + \epsilon_i$$

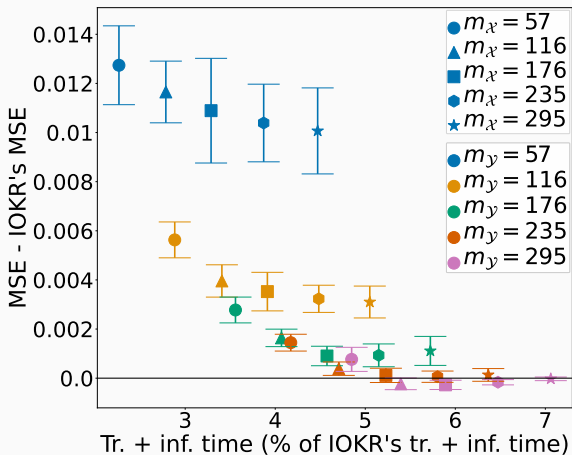
4)  $20/n$ -SR input and output sketches (sub-Gaussian)



# Synthetic least squares regression



# Synthetic least squares regression



# Multi-Label Classification: Statistical Performance

**Table 1:**  $F_1$  score on tag prediction from text data.

Method	Bibtex	Bookmarks
SISOKR	$44.1 \pm 0.07$	<b><math>39.3 \pm 0.61</math></b>
ISOKR	$44.8 \pm 0.01$	NA
SIOKR	$44.7 \pm 0.09$	$39.1 \pm 0.04$
IOKR	<b>44.9</b>	NA
LR	37.2	30.7
NN	38.9	33.8
SPEN	42.2	34.4
PRLR	44.2	34.9
DVN	44.7	37.1

# Multi-Label Classification: Computational Performance

**Table 2:** Comparison of training/inference computation times (in seconds).

Method	Bibtex	Bookmarks
SISOKR	<b><math>1.41 \pm 0.03</math> / <math>0.46 \pm 0.01</math></b>	<b><math>118 \pm 1.5</math> / <math>20 \pm 0.2</math></b>
ISOKR	$2.51 \pm 0.06$ / $0.58 \pm 0.01$	NA
SIOKR	$1.99 \pm 0.07$ / $1.22 \pm 0.03$	$354 \pm 2.1$ / $297 \pm 2.1$
IOKR	2.54 / 1.18	NA

## Synthetic and real-world experiments: take-home messages

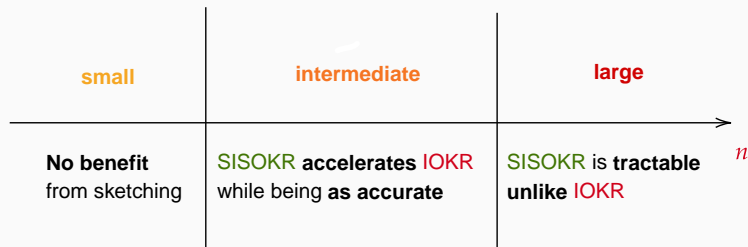
- 1) a) Input **sketching**: mainly accelerates the **training** phase
- 1) b) Output **sketching**: accelerates the **inference** phase

## Synthetic and real-world experiments: take-home messages

- 1) a) Input **sketching**: mainly accelerates the **training** phase
- 1) b) Output **sketching**: accelerates the **inference** phase
- 2) Optimal computational/statistical trade-off: statistical performance **converges** when  $m_x/m_y$  increases  $\implies$  **no need to set them too high!**

# Synthetic and real-world experiments: take-home messages

- 1) a) Input **sketching**: mainly accelerates the **training** phase
- 1) b) Output **sketching**: accelerates the **inference** phase
- 2) Optimal computational/statistical trade-off: statistical performance **converges** when  $m_x/m_y$  increases  $\implies$  no need to set them too high!
- 3) Benefits from **sketching** w.r.t. the **number of training data  $n$** :



## Conclusion

---



# Conclusion

- SISOKR: sketching on both input/output kernels to accelerate both training/inference steps
- Sketching as a way to build orthogonal projectors onto low-dimensional subspace of the feature space
- Excess risk bound leading to a consistent theoretical analysis of SISOKR
- Experiments: SISOKR accelerates IOKR or make it tractable

## References

---

- Brouard, C., d'Alché-Buc, F., and Szafranski, M. (2011). Semi-supervised penalized output kernel regression for link prediction. In *International Conference on Machine Learning (ICML)*, pages 593–600.
- Brouard, C., Shen, H., Dührkop, K., d'Alché-Buc, F., Böcker, S., and Rousu, J. (2016a). Fast metabolite identification with input output kernel regression. *Bioinformatics*, 32(12):28–36.
- Brouard, C., Szafranski, M., and D'Alché-Buc, F. (2016b). Input output kernel regression: supervised and semi-supervised structured output prediction with operator-valued kernels. *The Journal of Machine Learning Research*, 17(1):6105–6152.

- Caldarelli, E., Chatalic, A., Colomé, A., Molinari, C., Ocampo-Martinez, C., Torras, C., and Rosasco, L. (2024). Linear quadratic control of nonlinear systems with koopman operator learning and the nyström method.
- Ciliberto, C., Rosasco, L., and Rudi, A. (2020). A general framework for consistent structured prediction with implicit loss embeddings. *J. Mach. Learn. Res.*, 21(98):1–67.
- Cortes, C., Mohri, M., and Weston, J. (2005). A general regression technique for learning transductions. In *International Conference on Machine Learning (ICML)*, pages 153–160.

## References iii

- El Ahmad, T., Brogat-Motte, L., Laforgue, P., and d'Alché Buc, F. (2024). Sketch in, sketch out: Accelerating both learning and inference for structured prediction with kernels. In Dasgupta, S., Mandt, S., and Li, Y., editors, *Proceedings of The 27th International Conference on Artificial Intelligence and Statistics*, volume 238 of *Proceedings of Machine Learning Research*, pages 109–117. PMLR.
- Kadri, H., Ghavamzadeh, M., and Preux, P. (2013). A generalized kernel approach to structured output learning. In Dasgupta, S. and McAllester, D., editors, *Proceedings of the 30th International Conference on Machine Learning*, volume 28 of *Proceedings of Machine Learning Research*, pages 471–479, Atlanta, Georgia, USA. PMLR.
- Katakis, I., Tsoumakas, G., and Vlahavas, I. (2008). Multilabel text classification for automated tag suggestion. In *Proceedings of the ECML/PKDD*, volume 18, page 5. Citeseer.

- Korba, A., Garcia, A., and d'Alché-Buc, F. (2018). A structured prediction approach for label ranking. In Bengio, S., Wallach, H., Larochelle, H., Grauman, K., Cesa-Bianchi, N., and Garnett, R., editors, *Advances in Neural Information Processing Systems*, volume 31. Curran Associates, Inc.
- Meanti, G., Chatalic, A., Kostic, V. R., Novelli, P., massimiliano pontil, and Rosasco, L. (2023). Estimating koopman operators with sketching to provably learn large scale dynamical systems. In *Thirty-seventh Conference on Neural Information Processing Systems*.
- Rahimi, A. and Recht, B. (2007). Random features for large scale kernel machines. *NIPS*, 20:1177–1184.

## References v

- Rudi, A., Camoriano, R., and Rosasco, L. (2015). Less is more: Nyström computational regularization. *Advances in Neural Information Processing Systems*, 28.
- Schymanski, E., Ruttkies, C., and Krauss, M. e. a. (2017). Critical assessment of small molecule identification 2016: automated methods. *J Cheminform*, 9:22.
- Sriperumbudur, B. K. and Szabó, Z. (2015). Optimal rates for random fourier features. In *Proceedings of the 28th International Conference on Neural Information Processing Systems - Volume 1, NIPS'15*, page 1144–1152, Cambridge, MA, USA. MIT Press.
- Weston, J., Chapelle, O., Vapnik, V., Elisseeff, A., and Schölkopf, B. (2003). Kernel dependency estimation. In Becker, S., Thrun, S., and Obermayer, K., editors, *Advances in Neural Information Processing Systems 15*, pages 897–904. MIT Press.

- Williams, C. and Seeger, M. (2001). Using the nyström method to speed up kernel machines. In Leen, T., Dietterich, T., and Tresp, V., editors, *Advances in Neural Information Processing Systems*, volume 13, pages 682–688. MIT Press.
- Yang, T., Li, Y.-f., Mahdavi, M., Jin, R., and Zhou, Z.-H. (2012). Nyström method vs random fourier features: A theoretical and empirical comparison. In Pereira, F., Burges, C. J. C., Bottou, L., and Weinberger, K. Q., editors, *Advances in Neural Information Processing Systems*, volume 25. Curran Associates, Inc.
- Yang, Y., Pilanci, M., Wainwright, M. J., et al. (2017). Randomized sketches for kernels: Fast and optimal nonparametric regression. *The Annals of Statistics*, 45(3):991–1023.

# Complexity of IOKR and SISOKR for various types of sketching

**Table 3:** Time and space complexities at training and inference for the IOKR and SISOKR algorithms with sub-sampling,  $p$ -sparsified ( $p \in (0, 1]$ ) or Gaussian sketching, for a test set of size  $n_{te}$  and a candidate set of size  $n_c$ , such that  $n_{te} \leq m_x, m_y < n \leq n_c$ . For the sake of simplicity, we omit the  $\mathcal{O}(\cdot)$  in the following.

Method	Training		Inference	
	Time	Space	Time	Space
IOKR	$n^3$	$n^2$	$n_{te}nn_c$	$nn_c$
SISOKR (sub-sampling)	$\max(m_x, m_y)n$	$\max(m_x, m_y)n$	$n_{te}m_y n_c$	$m_y n_c$
SISOKR ( $p$ -sparsified)	$\max(m_x, m_y)^2 pn$	$\max(m_x, m_y)pn$	$\max(n_{te}, nm_y p)m_y n_c$	$npm_y n_c$
SISOKR (Gaussian)	$\max(m_x, m_y)n^2$	$n^2$	$nm_y n_c$	$nn_c$



# Sketching sizes selection strategy

**Goal:** set the minimal value of  $m_{\mathcal{Z}}$  s.t. it captures the information contained in the empirical covariance operator

$$\widehat{C}_{\mathcal{Z}} = \frac{1}{n} \sum_{i=1}^n \psi_{\mathcal{Z}}(z_i) \otimes \psi_{\mathcal{Z}}(z_i)$$

**However:** computing the SVD of  $\widehat{C}_{\mathcal{Z}}$  is costing, i.e.  $\mathcal{O}(n^3)$  in time.

1. Approximate leverage scores of  $\widehat{C}_{\mathcal{X}}$  and  $\widehat{C}_{\mathcal{Y}}$
2. Empirical approach: given training/inference budgets of time  $T_{\text{tr}}/T_{\text{inf}}$ , set low  $m_{\mathcal{X}}$  and  $m_{\mathcal{Y}}$  and evaluate the performance of  $\tilde{f}$  until reaching one of the following condition:
  - convergence of the performance of  $\tilde{f}$
  - training time attains  $T_{\text{tr}}$  or inference time attains  $T_{\text{te}}$

## Selection of $m_{\mathcal{X}}$

$\tilde{h}^{\text{SIOKR}}(x) = \sum_{i=1}^n \tilde{\alpha}_i^{\text{SIOKR}}(x) \psi_{\mathcal{Y}}(y_i)$  where

$$\tilde{\alpha}^{\text{SIOKR}}(x) = K_{\mathcal{X}} R_{\mathcal{X}}^{\top} (R_{\mathcal{X}} K_{\mathcal{X}}^2 R_{\mathcal{X}}^{\top} + n \lambda R_{\mathcal{X}} K_{\mathcal{X}} R_{\mathcal{X}}^{\top})^{\dagger}$$

Set the optimal  $m_{\mathcal{X}}$  according to a training budget of time  $T_{\text{tr}}$  and the performance of  $\tilde{h}^{\text{SIOKR}}$  in terms of surrogate regression error on the validation set, i.e. minimizing

$$\begin{aligned} & \sum_{i=1}^{n_{\text{val}}} \left\| \tilde{h}^{\text{SIOKR}}(x_i^{\text{val}}) - \psi_{\mathcal{Y}}(y_i^{\text{val}}) \right\|_{\mathcal{H}_{\mathcal{Y}}}^2 \\ &= \sum_{i=1}^{n_{\text{val}}} \tilde{\alpha}^{\text{SIOKR}}(x_i^{\text{val}})^{\top} K_{\mathcal{Y}} \tilde{\alpha}^{\text{SIOKR}}(x_i^{\text{val}}) - 2 \tilde{\alpha}^{\text{SIOKR}}(x_i^{\text{val}})^{\top} k_{\mathcal{Y}}^{y_i^{\text{val}}} + k_{\mathcal{Y}}(y_i^{\text{val}}, y_i^{\text{val}}) \end{aligned}$$

$\implies$  allows to cope with the inference phase

## Selection of $m_y$

Set the optimal  $m_y$  according to an inference budget of time  $T_{\text{inf}}$  and the performance of the *perfect*  $h$  estimator on the validation set, i.e.

$$h : (x, y) \mapsto \tilde{P}_Y \psi_Y(y)$$

$$f(x_i^{\text{val}}) = y_j^c \quad \text{where} \quad j = \arg \max_{1 \leq j \leq n_c} [K_Y^{\text{val, tr}} R_Y^T \tilde{K}_Y^\dagger R_Y K_Y^{\text{tr, c}}]_{ij}$$

$\implies$  allows to cope with the training phase

# Theory: previous works and differences

Rudi et al. (2015):

1. **scalar** kernel Ridge regression
2. sketching **only** applied in the **input** feature space
3. **Nyström** approximation with **uniform** or **approximate leverage scores** sampling

Ciliberto et al. (2020):

1. **vector-valued** kernel Ridge regression, with possibly infinite-dimensional outputs
2. **no approximation** considered

This work (El Ahmad et al., 2024):

1. **vector-valued** kernel Ridge regression, with possibly infinite-dimensional outputs
2. sketching applied in **both** the **input and output** feature space
3. generic **sub-Gaussian** sketches

Related recent works on Koopman operators: (Meanti et al., 2023; Caldarelli et al., 2024)

# SISOKR excess risk bound

Theorem (El Ahmad et al., 2024)

Let  $\delta \in [0, 1]$ ,  $n \in \mathbb{N}$  sufficiently large such that  $\lambda = n^{-1/(1+\gamma_x)} \geq \frac{9\kappa_x^2}{n} \log(\frac{n}{\delta})$ . Under **Asm. 1, 2, 3 and 4**, the following holds with probability at least  $1 - \delta$

$$\mathbb{E}[\|\tilde{h}(x) - h^*(x)\|_{\mathcal{H}_y}^2]^{\frac{1}{2}} \leq S(n) + c_2 A_{\rho_x}^{\psi_x}(\tilde{P}_X) + A_{\rho_y}^{\psi_y}(\tilde{P}_Y)$$

where

$$S(n) = c_1 \log(4/\delta) n^{-\frac{1}{2(1+\gamma_x)}} \quad (\text{regression error})$$

$$A_{\rho_z}^{\psi_z}(\tilde{P}_Z) = \mathbb{E}_Z[\|(\tilde{P}_Z - I_{\mathcal{H}_Z})\psi_Z(Z)\|_{\mathcal{H}_Z}^2]^{\frac{1}{2}} \quad (\text{sketching reconstruction error})$$

and  $c_1, c_2 > 0$  are constants independent of  $n$  and  $\delta$  defined in the proofs.

# Sub-Gaussian sketching reconstruction error

Theorem (El Ahmad et al., 2024)

Under **Asm. 1, 2, 3 and 4**, for  $\delta \in (0, 1/e]$ ,  $n \in \mathbb{N}$  sufficiently large such that  $\frac{9}{n} \log(n/\delta) \leq n^{-\frac{1}{1+\gamma_Z}} \leq \|C_Z\|_{\text{op}}/2$ , then if

$$m_Z \geq c_4 \max \left( \nu_Z^2 n^{\frac{\gamma_Z + \mu_Z}{1+\gamma_Z}}, \nu_Z^4 \log(1/\delta) \right),$$

then with probability  $1 - \delta$

$$\mathbb{E}_Z \left[ \left\| (\tilde{P}_Z - I_{\mathcal{H}_Z}) \psi_Z(Z) \right\|_{\mathcal{H}_Z}^2 \right] \leq c_3 n^{-\frac{1-\gamma_Z}{1+\gamma_Z}}$$

where  $c_3, c_4 > 0$  are constants independents of  $n, m_Z, \delta$  defined in the proofs.

# Multi-label classification

**Bibtex** and **Bookmarks** (Katakis et al., 2008): tag recommendation problems

**Mediamill**: detection of semantic concepts in a video

**Table 4:** Multi-label data sets description.

Data set	$n$	$n_{te}$	$n_{features}$	$n_{labels}$
Bibtex	4 880	2 515	1 836	159
Bookmarks	60 000	27 856	2 150	298
Mediamill	30 993	12 914	120	101

# Multi-label classification: statistical performance

**Table 5:**  $F_1$  scores on tag prediction from text data.

Method	Bibtex	Bookmarks	Mediamill
LR	37.2	30.7	NA
SPEN	42.2	34.4	NA
PRLR	44.2	34.9	NA
DVN	44.7	37.1	NA
SISOKR	$44.1 \pm 0.07$	<b><math>39.3 \pm 0.61</math></b>	$57.26 \pm 0.04$
ISOKR	$44.8 \pm 0.01$	NA	$58.02 \pm 0.01$
SIOKR	$44.7 \pm 0.09$	$39.1 \pm 0.04$	$57.33 \pm 0.04$
IOKR	<b>44.9</b>	NA	<b>58.17</b>



# Multi-label classification: computational performance

**Table 6:** Training/inference times (in seconds).

Method	Bibtex	Bookmarks	Mediamill
SISOKR	<b><math>1.41 \pm 0.03</math> / <math>0.46 \pm 0.01</math></b>	<b><math>118 \pm 1.5</math> / <math>20 \pm 0.2</math></b>	<b><math>66 \pm 0.1</math> / <math>4 \pm 0.01</math></b>
ISOKR	$2.51 \pm 0.06$ / $0.58 \pm 0.01$	NA	$636 \pm 3.7$ $9 \pm 0.2$
SIOKR	$1.99 \pm 0.07$ / $1.22 \pm 0.03$	$354 \pm 2.1$ / $297 \pm 2.1$	$199 \pm 0.1$ / $121 \pm 0.02$
IOKR	$2.54$ / $1.18$	NA	$621$ / $204$

# Metabolite identification

**Inputs:** tandem mass spectra of metabolites

**Outputs:** molecular structures, i.e. fingerprints, encoded by binary vectors of length  $d = 7593 \rightarrow$  **probability product kernel**

$n = 5579$  and each molecule is associated with a specific candidate set: **median size = 292** and **largest = 36918** fingerprints  $\rightarrow$

**Gaussian-Tanimoto kernel**

Method	kernel loss	Top-1   5   10 accuracies	training	inference
SPEN	$0.537 \pm 0.008$	25.9%   54.1%   64.3%	NA	NA
SISOKR	$0.566 \pm 0.007$	25.1%   54.2%   64.7%	$4.05 \pm 0.05$	<b><math>1112 \pm 29</math></b>
ISOKR	$0.509 \pm 0.009$	28.0%   58.9%   68.9%	$6.25 \pm 50.31$	$1133 \pm 32$
SIOKR	$0.492 \pm 0.008$	29.5%   61.3%   70.9%	<b><math>1.25 \pm 0.02</math></b>	$1179 \pm 37$
IOKR	<b><math>0.486 \pm 0.008</math></b>	<b>29.6%   61.6%   71.4%</b>	$3.54 \pm 0.15$	$1191 \pm 38$