

Neural Inverse Knitting: From Images to Manufacturing Instructions

Alexandre Kaspar*, Tae-Hyun Oh*, Liane Makatura, Petr Kellnhofer and Wojciech Matusik













Example: Sweater [Ministry of Supply]



Example: Scarf [Kniterate]





Industrial Knitting

• Whole garments from scratch



Industrial Knitting

- Control of individual needles
- Whole garments from scratch

Knitted Garment & Patterns

Many garments are knitted:

- Beanies, scarves
- Gloves, socks and underwear
- Sweaters, sweatpants

Current machines can create those garments **seamlessly** (no sewing needed).























Knitted Garment & Patterns

Those garments have **various types** of surface **patterns** (knitting patterns).

These can be fully controlled by industrial knitting machine.

= User customization!



Machine Knitting Programming

Low-level machine code requires skilled experts

= Knitting masters



Scenario

1. User takes picture of knitting pattern



Scenario

- 1. User takes picture of knitting pattern
- 2. System creates knitting instructions



Scenario

- 1. User takes picture of knitting pattern
- 2. System creates knitting instructions
- 3. User reuses pattern for new garment



Machine Knitting

Brief background

Machine Knitting Terminology



Illustration from [Underwood09]

Machine Knitting Terminology



Illustration from [Underwood09]

Machine Knitting Terminology

V-bed machine & knitting bed





"Tuck" operation

Illustration

from [Narayanan18]



"Knit"

operation

Illustration

from [Narayanan18]

"Transfer" operation



Illustration from [Narayanan18]

Racking = offsetting between the two beds



The Data & its Acquisition

For 2D machine knitting programs

2D Knitting Pattern Programs

Image of 20x20 pattern



20x20 pattern program



| Needles \rightarrow

"Pixels" are per-needle instructions over time

Knitting Pattern DSL

Domain Specific Language (DSL) for regular knitting patterns



Knitting Pattern DSL



DSL: from regular grids to sequences

Full rows of operations are executed at once with the following sequence:

- 1. Move "current stitches" to the operation side (front | back)
- 2. Apply "needle operation" (knit | tuck | miss)
- 3. Transfer moving stitches to back bed (cross | move | stack)
- 4. Apply *sequence* of moves depending on the operations (cross | move)
- 5. Bring back all stitches to front bed (purl | cross | move | stack)

DSL: from regular grids to sequences

Full rows of operations are executed at once with the following sequence:

Encoded by operation type:

- 1. Move "current stitches" to the opera **Move** = relative order not important
 - Cross = relative order defined by
 Group and "order" (upper blower)
- Apply "needle operation" (knit | tuckgroup and "order" (upper | lower)
- 3. Transfer moving stitches to back bed (cross | move | stack)
- 4. Apply *sequence* of operation-related moves (cross | move)
- 5. Bring back all stitches to front bed (purl | cross | move | stack)

Dataset: Initial Attempt

Individual 20x20 patterns

High-quality registration

- From color frame
- Still not per-stitch...

Total: ~200 patterns Time: ~1 month (intern) = not enough data!



Dataset: Better Attempt

Capture setup with steel rods to normalize tension





Dataset Content

- Paired instructions with real (2,088) and synthetic (14,440) images.
- Synthetic data from automatic screen capture of KnitPaint (Shima's software)



Machine Learning Details

Using two different types of supervision data

Learning Problem

Mapping **images** to discrete **instruction maps**

= CE loss minimization

Using two domains of data (one real, one synthetic)

= How to best combine both

With probability at least $1-\delta$

$$\frac{1}{2} \frac{|\mathcal{L}_T(\hat{h}, y) - \mathcal{L}_T(h_T^*, y)|}{\text{Generalization gap}} \leq \alpha \left(\text{disc}_{\mathcal{H}}(\mathcal{D}_S, \mathcal{D}_T) + \lambda \right) + \epsilon$$

$$\text{Ideal min.}$$

With probability at least $1-\delta$

$$\frac{\frac{1}{2}|\mathcal{L}_{T}(\hat{h}, y) - \mathcal{L}_{T}(h_{T}^{*}, y)|}{\text{Generalization gap}} \leq \alpha \left(\text{disc}_{\mathcal{H}}(\mathcal{D}_{S}, \mathcal{D}_{T}) + \lambda \right) + \epsilon$$

$$\text{Ideal min.}$$

Empirical min. $\arg \min_{h} \alpha \mathcal{L}_{\hat{S}}(h, y) + (1 - \alpha) \mathcal{L}_{\hat{T}}(h, y)$

With probability at least $1-\delta$

$$\frac{1}{2} |\mathcal{L}_T(\hat{h}, y) - \mathcal{L}_T(h_T^*, y)| \\ \leq \alpha \left(\operatorname{disc}_{\mathcal{H}}(\mathcal{D}_S, \mathcal{D}_T) + \lambda \right) + \epsilon$$

With probability at least $1-\delta$

$$\frac{1}{2} |\mathcal{L}_T(\hat{h}, y) - \mathcal{L}_T(h_T^*, y)| \le \alpha \left(\operatorname{disc}_{\mathcal{H}}(\mathcal{D}_S, \mathcal{D}_T) + \lambda \right) + \epsilon$$

$$\epsilon(m, \alpha, \beta, \delta) = \sqrt{\frac{1}{2m} \left(\frac{\alpha^2}{\beta} + \frac{(1-\alpha)^2}{1-\beta}\right) \log(\frac{2}{\delta})},$$

Hyper-parameter dependent term
Generalization Bound with Two Domains

With probability at least $1-\delta$

$$\frac{1}{2} |\mathcal{L}_T(\hat{h}, y) - \mathcal{L}_T(h_T^*, y)| \\ \leq \alpha \left(\operatorname{disc}_{\mathcal{H}}(\mathcal{D}_S, \mathcal{D}_T) + \lambda \right) + \epsilon$$
$$\lambda = \min_{h \in \mathcal{H}} \mathcal{L}_S(h, y) + \mathcal{L}_T(h, y).$$

Ideal error of the combined losses

Generalization Bound with Two Domains

With probability at least $1 - \delta$

$$\frac{1}{2} |\mathcal{L}_T(\hat{h}, y) - \mathcal{L}_T(h_T^*, y)| \leq \alpha \left(\operatorname{disc}_{\mathcal{H}}(\mathcal{D}_S, \mathcal{D}_T) + \lambda \right) + \epsilon$$

Discrepancy between distributions

$$\operatorname{disc}_{\mathcal{H}}(\mathcal{D}_{S}, \mathcal{D}_{T}) = \max_{h, h' \in \mathcal{H}} |\mathcal{L}_{\mathcal{D}_{S}}(h, h') - \mathcal{L}_{\mathcal{D}_{T}}(h, h')|$$

Data distributions

• Two different distribution types

Real data



Data distributions

• Two different distribution types



Real data

Synthetic data











Network composition



Results

Qualitative and quantitative evaluation







										75				λÐ	IA1	8.9	88	ы	Αī	6	aa	61	ATA
		6								35				80	a	80	80	IA1	81	â	aa	61	A18
					X									\mathbf{M}	a	89	80	IA1	81	â	10	63	818
		75								755				80	Ø.	80	80	IA1	81	â	80	61	A14
		75								75				80	i Al	av	80	IA1	81	à	80	o_1	114
		75								755				\mathbf{n}	M	\mathbf{a}	m	m	81	ä	30)	882	32
		75	X	*			X	X		7.5				50	n	an a	m	m	61	ð	-99	2	$\mathbf{\Omega}$
		7.5	0				0			7.5				875	n	av	m	n.	64	A	30	15	89
		7.5								7.5				823	m	av,	82	63	24	×	30	69	242
		7.5								7.5				874	en e	w	80	63	24	X	88	19	80
			0				Ň							822	03	W	80	624	29	A	29	20	-09
			0				0			75				823	22	QV)	80	624	69	A	193	10.	40
		75								75				822	64	N	M	604	ω.	2	878	25	640
		0								0				80	64	ΛV	NA	604	ω.	æ	26	2.5	842
			19	62			17	3.2						10	64	ΛV	86	201	λ.		49	30	616
			×	X			X	X						80	64	M	86	161	λ1	ଳ	60	40	816
														16	64	λV	86	161	λ.	F	10	20	88
														30	ы	λV	88	141	λŪ	ଳ	90	81	814
			Ω				Ω	Ω		0				30	ы	λ.	44	ы	81	କ	-83	61	816
			X	Ω			X	X	\times					343	ы	λ.	80	ы	81	କ	123	61	λ14
										25				MA	IAI	λ1/	10	IÅI	81	G	19.	0	A14
	0	-05	0	3	U	3	0	25		3		3		22	20	sa)	2	20	Sa:	a	82	20	2
														\sim	2	5	α	20	10	â	24	20	-
5	3	75	3	75	3	3	3	3	3	3	3	3	705	\sim	2	8	2	\sim	20	۵	\sim	20	2
5	5	3	75	6	3	3	75	3	3	3	3	3	3	\simeq	2	3	à	22	1	Ó	$\dot{\sim}$	22	1
5	3	-05	705	75	3	75	0	75	3	75	75	75	3	\sim	22	ð	à	22	2	Ċ.	24	-	
5	75	-05	12	75	~	75	12	75	~	75	75	75	3	\sim	22	1	$\Delta \mathbf{Y}$	20	1	ð	20	22	~
5	75	75	ω	75	-05	7.5	ω	75	7.5	7.5	7.5	7.5	75	\sim	2	3	4	2	2	à	22	10	5.5
	7.5	7.5	75	7.5	7.5		7.5	w.	7.5	75	7.5	75	-ŪS	\sim	2	3			2	ð	\sim	5	7
5	7.5	7.5	745	7.5	745	7-5	7.5	7.5	7.5	7.5	7.5	7.5	7.5	24	4	NA.	2	22	2	A	\sim	1	2
٠F	ŭ	7.5	7.5	7.5		7.5		7.5	7.5	7.5	7.5	7.5	7.5	70	44	ŝ	\sim	100	а.	8	~		
5	7.5	75	75	7.5	7-5	7.5	7.5	7.5	7.5	7.5	7.5	7.5	7.5	~	12	3	24	~	10	2	~	1	20
	7.5		7-5	Ň	7.5		7.5	7.5	7.5	7.5	7.5	7.5	7.5	20	10	3	20	- 7	28	2	23	22	8.0
5	7.5	7.5	75	7.5		7.5	75	75	75	7.5	75	75		70	424	3	20	29	2	2	20	1.1	22
5	č	7.5	ŭ	7.5	ŭ	7.5	75	75	Ň	7.5	ň	7.5		え	10	3	*	22	3	52		20	100
5	7.5		75	75	75	75	75	75	7.5	75	7.5	7.5	7.5	10	10	3	20	10	37	83	22	22	370
5	7.5	č	25	-	7.5	75	75		25	Ň	7.5	7.5	7.5	\mathcal{T}	10	33	*	22	32		20	×.,	88
5	25	75	75	75	25	75	75	75	25	75	25	75	75	20	10	27	20	10	32	50	20	20	323
5	75	75	75	75	75	75	75	75	75	75	75	75		20	2	23	20	22	22	2	20	25	22
5	0	-0-	0	-0-	-0-	-0-	0	-0-	-0-	20	-0-	-0-	0	23	-34	58	22	10	5	2	20	-32	88
2	0	0	0	0	0	0	0	0	0	0	0	0	0	20	20	5	20	202	82	9	20	- 1	2
9	0	0	\odot	0	0	0	\odot	0	0	0	0	0	0	24	1	R.	94	- 6	6	R	24	20	50
		12												36	16	37	86	99	69	10	89	65	690
		52				52				52				30	м	M	18	633	Ш		16	6.0	11
														30	ы	λV	db	88.	60	64	29	88	660
		Ω				52				52				30	64	80		83	лı		14	6.3	8.8
														363	61	λ1/	db	88.	80	68	810	-51	hab
														310	ы	a	а.	636	ðТ	2	181	631	2.8
														303	Ø.	80	66	88.	81	1.1	20	193	8.60
														14	A.	1.1		16	1		14		1
		2				1				1				14	A.	11	180	10	63	11	14	10	10
												~		10	IA.	11		10.			23		
		1				1								10	A.	11	18		83	1	10	8	1.1
														10	A.	11		14		1	12		
		1				1.				1.				10		11		10	13	1	10		1.1
														10	10	1	100	10		1	12	1	
		-				-				1				10		1	1	10		1	10		1
														10		1	1	12		1	p)	51	1
		-				N				N				10	1	1	1	0		1	18	10	1
														10	1	1	1	12			Pa	5.	
		1				1				N				10	1	2	1	01	14	1	10	1	1
														10	1	1	19	12	1	14	100	1	14
																and the second s			and the second s	_	and the second sec	-	



Architecture variations



	Mathad	Accura	acy (%)	Perceptual		
	Wethod	Full	FG	SSIM	PSNR [dB]	
(a1)	CycleGAN (Zhu et al., 2017)	57.27	24.10	0.670	15.87	
(a2)	Pix2Pix (Isola et al., 2017)	56.20	47.98	0.660	15.95	
(a3)	UNet (Ronneberger et al., 2015)	89.65	63.99	0.847	21.21	
(a4)	Scene Parsing (Zhou et al., 2018)	91.58	73.95	0.876	22.64	
(a5)	S+U (Shrivastava et al., 2017)	91.32	71.00	0.864	21.42	
(b1)	Img2prog (real only) with CE	91.57	71.37	0.866	21.62	
(b2)	Img2prog (real only) with MILCE	91.74	72.30	0.871	21.58	
(c1)	Refiner + Img2prog ($\alpha = 0.9$)	93.48	78.53	0.894	23.28	
(c2)	Refiner + Img2prog ($lpha=2/3$)	93.58	78.57	0.892	23.27	
(c3)	Refiner + Img2prog ($\alpha = 0.5$)	93.57	78.30	0.895	23.24	
(c4)	Refiner + Img2prog ($lpha=1/3$)	93.19	77.80	0.888	22.72	
(c5)	Refiner + Img2prog ($\alpha = 0.1$)	92.42	74.15	0.881	22.27	
(d1)	Refiner + Img2prog++ ($\alpha = 0.5$)	94.01	80.30	0.899	23.56	

	Method	Accura	acy (%)	Pe	rceptual
	Wethou	Full	FG	SSIM	PSNR [dB]
(a1)	CycleGAN (Zhu et al., 2017)	57.27	24.10	0.670	15.87
(a2)	Pix2Pix (Isola et al., 2017)	56.20	47.98	0.660	15.95
(a3)	UNet (Ronneberger et al., 2015)	89.65	63.99	0.847	21.21
(a4)	Scene Parsing (Zhou et al., 2018)	91.58	73.95	0.876	22.64
(a5)	S+U (Shrivastava et al., 2017)	91.32	71.00	0.864	21.42
(b1)	Img2prog (real only) with CE	91.57	71.37	0.866	21.62
(b2)	Img2prog (real only) with MILCE	91.74	72.30	0.871	21.58
(c1)	Refiner + Img2prog ($\alpha = 0.9$)	93.48	78.53	0.894	23.28
(c2)	Refiner + Img2prog ($\alpha = 2/3$)	93.58	78.57	0.892	23.27
(c3)	Refiner + Img2prog ($\alpha = 0.5$)	93.57	78.30	0.895	23.24
(c4)	Refiner + Img2prog ($\alpha = 1/3$)	93.19	77.80	0.888	22.72
(c5)	Refiner + Img2prog ($\alpha = 0.1$)	92.42	74.15	0.881	22.27
(d1)	Refiner + Img2prog++ ($\alpha = 0.5$)	94.01	80.30	0.899	23.56

	Mathad	Accura	acy (%)	Perceptual		
	Method	Full	FG	SSIM	PSNR [dB]	
(a1)	CycleGAN (Zhu et al., 2017)	57.27	24.10	0.670	15.87	
(a2)	Pix2Pix (Isola et al., 2017)	56.20	47.98	0.660	15.95	
(a3)	UNet (Ronneberger et al., 2015)	89.65	63.99	0.847	21.21	
(a4)	Scene Parsing (Zhou et al., 2018)	91.58	73.95	0.876	22.64	
(a5)	S+U (Shrivastava et al., 2017)	91.32	71.00	0.864	21.42	
(b1)	Img2prog (real only) with CE	91.57	71.37	0.866	21.62	
(b2)	Img2prog (real only) with MILCE	91.74	72.30	0.871	21.58	
(c1)	Refiner + Img2prog ($\alpha = 0.9$)	93.48	78.53	0.894	23.28	
(c2)	Refiner + Img2prog ($\alpha = 2/3$)	93.58	78.57	0.892	23.27	
(c3)	Refiner + Img2prog ($\alpha = 0.5$)	93.57	78.30	0.895	23.24	
(c4)	Refiner + Img2prog ($\alpha = 1/3$)	93.19	77.80	0.888	22.72	
(c5)	Refiner + Img2prog ($\alpha = 0.1$)	92.42	74.15	0.881	22.27	
(d1)	Refiner + Img2prog++ ($\alpha = 0.5$)	94.01	80.30	0.899	23.56	

]	Mathad	Accura	acy (%)	Perceptual		
	Wethou	Full	FG	SSIM	PSNR [dB]	
(a1)	CycleGAN (Zhu et al., 2017)	57.27	24.10	0.670	15.87	
(a2)	Pix2Pix (Isola et al., 2017)	56.20	47.98	0.660	15.95	
(a3)	UNet (Ronneberger et al., 2015)	89.65	63.99	0.847	21.21	
(a4)	Scene Parsing (Zhou et al., 2018)	91.58	73.95	0.876	22.64	
(a5)	S+U (Shrivastava et al., 2017)	91.32	71.00	0.864	21.42	
(b1)	Img2prog (real only) with CE	91.57	71.37	0.866	21.62	
(b2)	Img2prog (real only) with MILCE	91.74	72.30	0.871	21.58	
(c1)	Refiner + Img2prog ($\alpha = 0.9$)	93.48	78.53	0.894	23.28	
(c2)	Refiner + Img2prog ($lpha=2/3$)	93.58	78.57	0.892	23.27	
(c3)	Refiner + Img2prog ($\alpha = 0.5$)	93.57	78.30	0.895	23.24	
(c4)	Refiner + Img2prog ($lpha=1/3$)	93.19	77.80	0.888	22.72	
(c5)	Refiner + Img2prog ($\alpha = 0.1$)	92.42	74.15	0.881	22.27	
(d1)	Refiner + Img2prog++ ($\alpha = 0.5$)	94.01	80.30	0.899	23.56	

How much data is enough data?



Limitations

And potential solutions

Issue of scale, stretch and orientation

We assume a specific scale, stretch (of 20x20 stitches) and a bottom-up orientation of stitch courses.

Options:

- Explicit model scale, stretch and orientation
 = makes training more complicated
- Separate selection (using measure of "confidence")
 - = take large-scale image, and try space of scales / stretches / rot.

Attempt at scale selection (successful)



Crop scale [px]

Input variety

We only used Tamm 2/30 acrylic yarn.

How do we scale to more data, and more varieties of it?

Options:

- Simulation: need fast simulation of yarn (hard, or slow), hopefully as a differentiable renderer (within the network)
- Online yarn images: unsupervised way? Cycle-consistency?
 Additional side/weaker/stronger task?

Modeling Hard Constraints

Currently, output may have invalid instruction combinations. Tried to use penalty on valid 1st order neighborhood, but little impact.

Questions:

- How do we model hard constraints with a neural network?
- Split translation into instruction "potentials" and then select the actual instructions (e.g., using known knittability constraints)?
- Can we infer the syntax constraints automatically?
 - Note: non-trivial to specify beyond first-order neighborhood unless enough data is available...

Result Details

The great, the good, the not so good, and the ugly

Details: Perfect cases



Details: Minor errors (no semantic issue)



Details: Larger errors (but knittable)



Initial program

Initial sample

Inferred program

Final sample

Details: Larger errors (but knittable)



Initial program

Initial sample

Inferred program

Final sample

Details: (few) catastrophic failures (only 2)



Past and Future Work

Where it came from, and where it is going



Recent: Knitting Skeletons - CAD for Knitting







Initial	Range of wales	of cours
Area select	Union	neighbor
img mask	tile mask	Filter with no

(a) Knit _(-)_	(b) Purl 75	(c) Tuck
		ARAAAAA AAAAAA ARAAAAAA ARAAAAAA ARAAAAAA

(e) Move

(d) Miss

(f) Cross 📉

Now: Knit Sketching - Sketches within CAD

Work with sketches

- Wale flow
- Connectivity
- Stitch density
- Layers (sketches)
- Layers (patterns)

Generate data for the CAD system.

(with some efficient parameterization)



Next: InverseKnit++

Use sketch input capability to learn to map full knitted "shapes" directly into low-level knitting programs.

- More instruction irregularities
- Issue of occlusion (two-sided shapes)
- Ambiguity between shape and patterns



http://deepknitting.csail.mit.edu



-
Dataset Details

Instruction distribution and accuracies

Dataset: instruction statistics



Figure 5. Instruction counts in descending order, for synthetic and real images. Note the logarithmic scale of the Y axis.

Per-Instruction Accuracies

Table 2. **Performance of** *Refined*+*Img2prog*++ **measured per instruction over the test set.** This shows that even though our instruction distribution has very large variations, our network is still capable of learning some representation for the least frequent instructions (3 orders of magnitude difference for FR2, FL2, BR2, BL2 compared to K and P).

Instruction	K	Р	Т	Μ	FR1	FR2	FL1	FL2	BR1	BR2	BL1	BL2	XR+	XR-	XL+	XL-	S
Accuracy [%]	96.52	96.64	74.63	66.65	77.16	100.00	74.20	83.33	68.73	27.27	69.94	22.73	60.15	62.33	60.81	62.11	25.85
Frequency [%]	44.39	47.72	0.41	1.49	1.16	0.01	1.23	0.01	1.22	0.02	1.40	0.02	0.22	0.18	0.19	0.22	0.12

Architecture Details

Neural Networks and Losses

Actual Loss Function



Actual Loss Function

Our combined loss is the weighted sum

$$\mathcal{L} = \lambda_{\rm CE} \mathcal{L}_{\rm CE} + \lambda_{\rm Perc} \mathcal{L}_{\rm Perc} + \lambda_{\rm GAN} \mathcal{L}_{\rm GAN}$$
(5)

where we used the weights: $\lambda_{CE} = 3$, $\lambda_{Perc} = 0.02/(128)^2$ and $\lambda_{GAN} = 0.2$. The losses \mathcal{L}_{Perc} and λ_{GAN} are measured on the output of Refiner, while the loss λ_{CE} is measured on Img2prog.

Refiner Network



Figure 10. The illustration of the Refiner network architecture, where S # N denotes the stride size of # N, IN_ReLU indicates the Instance normalization followed by ReLU, Resblk is the residual block that consists of ConvS1-ReLU-ConvS1 with short-cut connection (He et al., 2016), Upsample is the nearest neighbor upsampling with the factor $2 \times$, F is the output channel dimension. If not mentioned, the default parameters for all the convolutions are the stride size of 2, F = 64, and the 3×3 kernel size.

Theorem 1

About the Generalization Gap

Definition 1: Discrepancy [Mansour 09]

Definition 1 (Discrepancy (Mansour et al., 2009)). Let \mathcal{H} be a class of functions mapping from \mathcal{X} to \mathcal{Y} . The discrepancy between two distribution \mathcal{D}_1 and \mathcal{D}_2 over \mathcal{X} is defined as

$$\operatorname{disc}_{\mathcal{H}}(\mathcal{D}_1, \mathcal{D}_2) = \max_{h, h' \in \mathcal{H}} \left| \mathcal{L}_{\mathcal{D}_1}(h, h') - \mathcal{L}_{\mathcal{D}_2}(h, h') \right|.$$
(6)

The discrepancy is symmetric and satisfies the triangle inequality, regardless of any loss function. This can be used to compare distributions for general tasks even including regression.

Lemma 1. Let h be a hypothesis in class \mathcal{H} , and assume that \mathcal{L} is symmetric and obeys the triangle inequality. Then

 $\left|\mathcal{L}_{\alpha}(h, y) - \mathcal{L}_{T}(h, y)\right| \leq \alpha \left(\operatorname{disc}_{\mathcal{H}}(\mathcal{D}_{S}, \mathcal{D}_{T}) + \lambda\right), \quad (7)$

where $\lambda = \mathcal{L}_S(h^*, y) + \mathcal{L}_T(h^*, y)$, and the ideal joint hypothesis h^* is defined as $h^* = \arg \min_{h \in \mathcal{H}} \mathcal{L}_S(h, y) + \mathcal{L}_T(h, y)$.

Proof. The proof is based on the triangle inequality of \mathcal{L} , and the last inequality follows the definition of the discrepancy.

 $|\mathcal{L}_{\alpha}(h,y) - \mathcal{L}_{T}(h,y)|$ Substitute $\mathscr{L}_{\alpha} = \alpha \mathscr{L}_{S} + (1 - \alpha) \mathscr{L}_{T}$ $=\alpha |\mathcal{L}_S(h, y) - \mathcal{L}_T(h, y)|$ $=\alpha |\mathcal{L}_{S}(h, y) - \mathcal{L}_{S}(h^{*}, h) + \mathcal{L}_{S}(h^{*}, h)$ $-\mathcal{L}_T(h^*,h) + \mathcal{L}_T(h^*,h) - \mathcal{L}_T(h,y)$ $\leq \alpha ||\mathcal{L}_S(h, y) - \mathcal{L}_S(h^*, h)| +$ $|\mathcal{L}_{S}(h^{*},h) - \mathcal{L}_{T}(h^{*},h)| + |\mathcal{L}_{T}(h^{*},h) - \mathcal{L}_{T}(h,y)||$ $\leq \alpha \left| \mathcal{L}_S(h^*, y) + \left| \mathcal{L}_S(h^*, h) - \mathcal{L}_T(h^*, h) \right| + \mathcal{L}_T(h^*, y) \right|$ $\leq \alpha \left(\operatorname{disc}_{\mathcal{H}}(\mathcal{D}_S, \mathcal{D}_T) + \lambda \right).$ (8)

Proof. The proof is based on the triangle inequality of \mathcal{L} , and the last inequality follows the definition of the discrepancy.

 $|\mathcal{L}_{\alpha}(h, y) - \mathcal{L}_{T}(h, y)|$ $= \alpha |\mathcal{L}_S(h, y) - \mathcal{L}_T(h, y)|$ Introduce $(\mathscr{L}_{S}(h^{*},h) - \mathscr{L}_{S}(h^{*},h))$ $= \alpha \left| \mathcal{L}_{S}(h, y) - \mathcal{L}_{S}(h^{*}, h) + \mathcal{L}_{S}(h^{*}, h) \right|$ and $(\mathscr{L}_{\tau}(h^{*},h))$ - $\mathscr{L}_{\tau}(h^{*},h))$ $-\mathcal{L}_T(h^*,h)+\mathcal{L}_T(h^*,h)-\mathcal{L}_T(h,y)$ $\leq \alpha | |\mathcal{L}_S(h, y) - \mathcal{L}_S(h^*, h)| +$ $|\mathcal{L}_{S}(h^{*},h) - \mathcal{L}_{T}(h^{*},h)| + |\mathcal{L}_{T}(h^{*},h) - \mathcal{L}_{T}(h,y)||$ $\leq \alpha \left| \mathcal{L}_S(h^*, y) + \left| \mathcal{L}_S(h^*, h) - \mathcal{L}_T(h^*, h) \right| + \mathcal{L}_T(h^*, y) \right|$ $\leq \alpha \left(\operatorname{disc}_{\mathcal{H}}(\mathcal{D}_S, \mathcal{D}_T) + \lambda \right).$ (8)

Proof. The proof is based on the triangle inequality of \mathcal{L} , and the last inequality follows the definition of the discrepancy.

 $|\mathcal{L}_{\alpha}(h, y) - \mathcal{L}_{T}(h, y)|$ $=\alpha |\mathcal{L}_S(h, y) - \mathcal{L}_T(h, y)|$ $=\alpha |\mathcal{L}_S(h, y) - \mathcal{L}_S(h^*, h) + \mathcal{L}_S(h^*, h)$ $-\mathcal{L}_T(h^*,h) + \mathcal{L}_T(h^*,h) - \mathcal{L}_T(h,y)$ Apply triangular inequality $\leq \alpha | |\mathcal{L}_S(h, y) - \mathcal{L}_S(h^*, h)| +$ $|\mathcal{L}_{S}(h^{*},h) - \mathcal{L}_{T}(h^{*},h)| + |\mathcal{L}_{T}(h^{*},h) - \mathcal{L}_{T}(h,y)||$ $\leq \alpha \left| \mathcal{L}_S(h^*, y) + \left| \mathcal{L}_S(h^*, h) - \mathcal{L}_T(h^*, h) \right| + \mathcal{L}_T(h^*, y) \right|$ $\leq \alpha \left(\operatorname{disc}_{\mathcal{H}}(\mathcal{D}_{S}, \mathcal{D}_{T}) + \lambda \right).$ (8)

Proof. The proof is based on the triangle inequality of \mathcal{L} , and the last inequality follows the definition of the discrepancy.

$$\begin{aligned} |\mathcal{L}_{\alpha}(h,y) - \mathcal{L}_{T}(h,y)| \\ = \alpha |\mathcal{L}_{S}(h,y) - \mathcal{L}_{T}(h,y)| \\ = \alpha |\mathcal{L}_{S}(h,y) - \mathcal{L}_{S}(h^{*},h) + \mathcal{L}_{S}(h^{*},h) \\ - \mathcal{L}_{T}(h^{*},h) + \mathcal{L}_{T}(h^{*},h) - \mathcal{L}_{T}(h,y)| \\ \leq \alpha |[\mathcal{L}_{S}(h,y) - \mathcal{L}_{S}(h^{*},h)] + \\ |\mathcal{L}_{T}(h^{*},h) - \mathcal{L}_{T}(h^{*},h)| + |[\mathcal{L}_{T}(h^{*},h) - \mathcal{L}_{T}(h,y)|] \\ \leq \alpha |[\mathcal{L}_{S}(h^{*},y) + |[\mathcal{L}_{S}(h^{*},h) - \mathcal{L}_{T}(h^{*},h)] + [\mathcal{L}_{T}(h^{*},y)] \\ \leq \alpha (\operatorname{disc}_{\mathcal{H}}(\mathcal{D}_{S},\mathcal{D}_{T}) + \lambda). \end{aligned}$$
(8)

Proof. The proof is based on the triangle inequality of \mathcal{L} , and the last inequality follows the definition of the discrepancy.

 $|\mathcal{L}_{\alpha}(h, y) - \mathcal{L}_{T}(h, y)|$ $=\alpha |\mathcal{L}_S(h, y) - \mathcal{L}_T(h, y)|$ $=\alpha \left| \mathcal{L}_{S}(h, y) - \mathcal{L}_{S}(h^{*}, h) + \mathcal{L}_{S}(h^{*}, h) \right|$ $-\mathcal{L}_T(h^*,h)+\mathcal{L}_T(h^*,h)-\mathcal{L}_T(h,y)$ $\leq \alpha ||\mathcal{L}_S(h, y) - \mathcal{L}_S(h^*, h)| +$ $|\mathcal{L}_{S}(h^{*},h) - \mathcal{L}_{T}(h^{*},h)| + |\mathcal{L}_{T}(h^{*},h) - \mathcal{L}_{T}(h,y)||$ $\leq \alpha \left| \mathcal{L}_S(h^*, y) + \left| \mathcal{L}_S(h^*, h) - \mathcal{L}_T(h^*, h) \right| + \mathcal{L}_T(h^*, y) \right|$ Definition of discrepancy + storing rest in λ $\leq \alpha \left(\operatorname{disc}_{\mathcal{H}}(\mathcal{D}_S, \mathcal{D}_T) + \lambda \right).$ (8)

Lemma 2: from [Ben-David 10]

Lemma 2 ((Ben-David et al., 2010)). For a fixed hypothesis h, if a random labeled sample of size m is generated by drawing βm points from \mathcal{D}_S and $(1 - \beta)m$ points from \mathcal{D}_T , and labeling them according to y_S and y_T respectively, then for any $\delta \in (0, 1)$, with probability at least $1 - \delta$ (over the choice of the samples),

$$|\hat{\mathcal{L}}_{\alpha}(h,y) - \mathcal{L}_{\alpha}(h,y)| \le \epsilon(m,\alpha,\beta,\delta), \tag{9}$$

where
$$\epsilon(m, \alpha, \beta, \delta) = \sqrt{\frac{1}{2m} \left(\frac{\alpha^2}{\beta} + \frac{(1-\alpha)^2}{1-\beta}\right) \log(\frac{2}{\delta})}.$$

The detail function form of ϵ will be omitted for simplicity. We can fix m, α , β , and δ when the learning task is specified, then we can treat $\epsilon(\cdot)$ as a constant.

Theorem 1: Generalization Gap

Theorem 1. Let \mathcal{H} be a hypothesis class, and S be a labeled sample of size m generated by drawing βm samples from \mathcal{D}_S and $(1 - \beta)m$ samples from \mathcal{D}_T and labeling them according to the true label y. Suppose \mathcal{L} is symmetric and obeys the triangle inequality. Let $\hat{h} \in \mathcal{H}$ be the empirical minimizer of $\hat{h} = \arg \min_h \hat{\mathcal{L}}_\alpha(h, y)$ on S for a fixed $\alpha \in [0, 1]$, and $h_T^* = \arg \min_h \mathcal{L}_T(h, y)$ the target error minimizer. Then, for any $\delta \in (0, 1)$, with probability at least $1 - \delta$ (over the choice of the samples), we have

$$\frac{1}{2}|\mathcal{L}_{T}(\hat{h}, y) - \mathcal{L}_{T}(h_{T}^{*}, y)| \leq \alpha \left(\operatorname{disc}_{\mathcal{H}}(\mathcal{D}_{S}, \mathcal{D}_{T}) + \lambda\right) + \epsilon,$$
(10)
where $\epsilon(m, \alpha, \beta, \delta) = \sqrt{\frac{1}{2m} \left(\frac{\alpha^{2}}{\beta} + \frac{(1-\alpha)^{2}}{1-\beta}\right) \log(\frac{2}{\delta})}, and$

$$\lambda = \min_{h \in \mathcal{H}} \mathcal{L}_{S}(h, y) + \mathcal{L}_{T}(h, y).$$

Proof of Theorem 1

Proof. We use Lemmas 1 and 2 for the bound derivation with their associated assumptions.

$$\mathcal{L}_{T}(\hat{h}, y) \leq \mathcal{L}_{\alpha}(\hat{h}, y) + \alpha \left(\operatorname{disc}_{\mathcal{H}}(\mathcal{D}_{S}, \mathcal{D}_{T}) + \lambda \right), \qquad (11)$$

$$(\operatorname{By Lemma 1}) \leq \hat{\mathcal{L}}_{\alpha}(\hat{h}, y) + \alpha \left(\operatorname{disc}_{\mathcal{H}}(\mathcal{D}_{S}, \mathcal{D}_{T}) + \lambda \right) + \epsilon, \qquad (12)$$

$$(\operatorname{By Lemma 2}) \leq \hat{\mathcal{L}}_{\alpha}(h_{T}^{*}, y) + \alpha \left(\operatorname{disc}_{\mathcal{H}}(\mathcal{D}_{S}, \mathcal{D}_{T}) + \lambda \right) + \epsilon, \qquad (13)$$

$$(\hat{h} = \arg\min_{h \in \mathcal{H}} \hat{\mathcal{L}}_{\alpha}(h))$$

$$\leq \mathcal{L}_{\alpha}(h_{T}^{*}, y) + \alpha \left(\operatorname{disc}_{\mathcal{H}}(\mathcal{D}_{S}, \mathcal{D}_{T}) + \lambda \right) + 2\epsilon, \qquad (14)$$

$$(\operatorname{By Lemma 2})$$

$$\leq \mathcal{L}_{T}(h_{T}^{*}, y) + 2\alpha \left(\operatorname{disc}_{\mathcal{H}}(\mathcal{D}_{S}, \mathcal{D}_{T}) + \lambda \right) + 2\epsilon, \qquad (15)$$

$$(\operatorname{By Lemma 1})$$

Cost of "swapping" is at least ...

$$|\mathcal{L}_{\alpha}(h, y) - \mathcal{L}_{T}(h, y)| \leq \alpha \left(\operatorname{disc}_{\mathcal{H}}(\mathcal{D}_{S}, \mathcal{D}_{T}) + \lambda \right)$$

(12) Specific sampling: $|\hat{\mathcal{L}}_{\alpha}(h, y) - \mathcal{L}_{\alpha}(h, y)| \le \epsilon(m, \alpha, \beta, \delta)$

Because
$$\hat{\mathcal{L}}_{\alpha}(\hat{h},y) \leq \hat{\mathcal{L}}_{\alpha}(h_{T}^{*},y)$$

Specific "de-sampling": $|\hat{\mathcal{L}}_{\alpha}(h, y) - \mathcal{L}_{\alpha}(h, y)| \leq \epsilon(m, \alpha, \beta, \delta)$

 $|\mathcal{L}_{\alpha}(h, y) - \mathcal{L}_{T}(h, y)| \le \alpha (\operatorname{disc}_{\mathcal{H}}(\mathcal{D}_{S}, \mathcal{D}_{T}) + \lambda)$