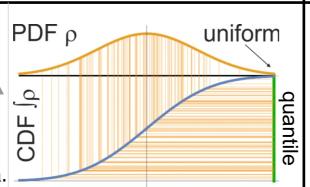
## Hierarchical Correlation Reconstruction Hierarchical Correlation Reconstruction HIERARN biology-inspired neural networks Logical two-way(2W): Maximum distribution and training ways Inspired neural networks of normalized variables Milli-directional propagation, of values/distributions, additional training ways Duda, JU arXiv:2405.05097 sides from hilli-directional propagation, of values/distributions, additional training ways Duda, JU arXiv:2405.05097 sides from hilli-directional propagation, of values/distributions, additional training ways Duda, JU arXiv:2405.05097 sides from hilli-directional propagation, of values/distributions, additional training ways parameters with neurons containing local joint distribution model as polynomial **Hierarchical Correlation Reconstruction** multi-directional propagation, of values/distrbutions, additional training ways

1) Biological neural networks (BNNs) seem qualitatively superior to current artificial (ANNs). Biological neurons have thousands of connections, also to itself, allowing to hide complex mathematics, also through bursting and cancellation of colliding action potentials. Instead of such complex spiking NN physical level, let us search for higher level "logical neuron" here: extracting "biologically plausible" mathematics representing neuron dynamics.

Parameters						
I ul ullicter 5	ANN <u>source</u>	BNN				
	input	dendrites	L: ALPOYERE			
Structure	weight	synapse	经投资政计论			
Structure	output	axon	LAK DEST			
	hidden layer	cell body (~86B)	The Course			
Loorning	very precise structures	they can tolerate				
Learning	and formatted data	ambiguity	Action At t=0ms Plasma membri			
	complex	simple				
Processor	high speed	low speed				
	one or a few	large number	L: Action t=3ms			
	separate from a	integrated	1: potential + + + +			
Memory	processor	into processor	Na'			
Memory	localized	distributed	K* t=7ms			
	non-content addressable	content-addressable	Action potential			
	centralized	distributed				
Computing	sequential	paralel	+ +			
	stored programs	self-learning	dendrites			
Reliability	very vulnerable	robust	nucleus			
Expertise	numerical and symbolic	perceptual	cell axon			
-	manipulations	nipulations problems <u>l:</u>	axon terminals			
Operating	well-defined	poorly defined	in			
Environment	well-constrained	un-constrained	$in_2 \longrightarrow \Sigma f$ out			
BNN learning $\neq$ ANN backpropagation $2W$ :						
	superiority: just "ta		bias two-			

4) For joint distribution data structure, it is very convenient to normalize variables to **nearly uniform in [0,1]** (e.g. by  $x \rightarrow CDF(x)$ or EDF by just sorting), then **represent their** joint density as just linear combination e.g. of polynomials, preferably in orthonormals basis allowing for independent MSE estimation.

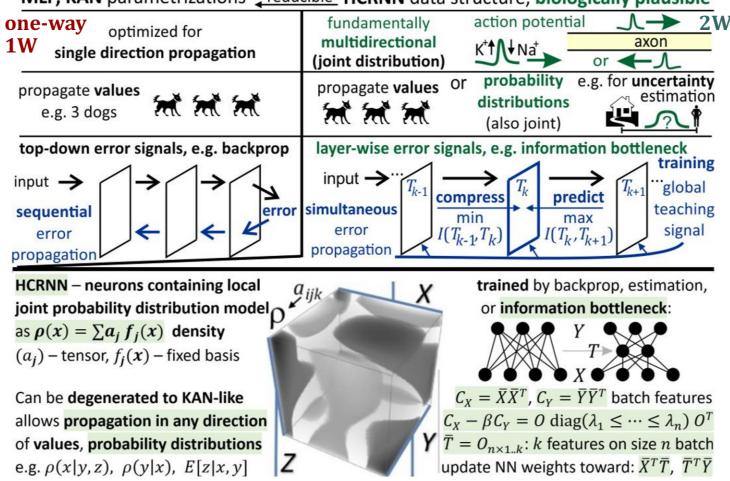


decompose dependencies into mixed moments – control, accuracy, credibilty

2) While standard logical neurons like MLP (Multi-Layer Perceptron), KAN are **one-way (1W)**: focused on **uni-directional propagation**, **biological axons** propagate in both directions (2W) e.g. for training, forming two-way neural **networks (2WNN)**. Also value prediction alone seems insufficient for animals it is crucial to estimate uncertainty, e.g. by propagating entire distributions. Additionally, backpropagation is viewed as not biologically plausible biology needs local layer-wise training ways like information bottleneck.

Spiking: physical level, but what is 'logical neuron': hidden mathematics?

MLP, KAN parametrizations \_ reducible HCRNN data structure, biologically plausible

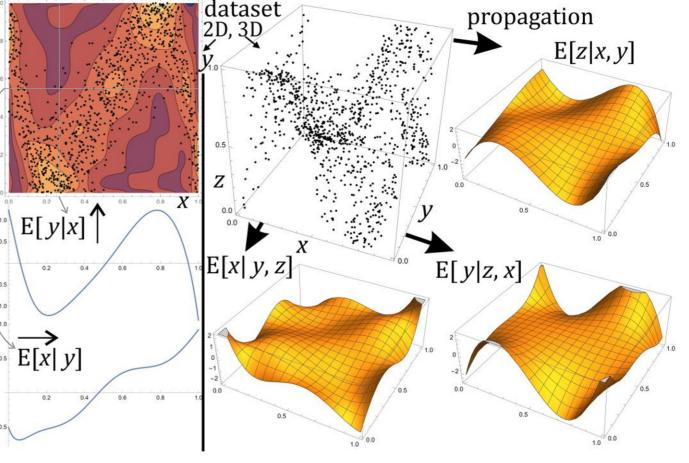


5) While local basis like Gaussians is useful for clustered data, for neurons there is sought e.g. XOR operation - trivial to realize with polynomial. Generally, global basis like polynomials usually gives much better cross-validation like below - generalizes better: searching for general data description like moments here, instead of just assuming that "new points will be close to old points" in local e.g. KDE or mixture models, also less convenient to represent.

**moments generalize well** log-likelihood: mean  $lg(\rho)$  on random 25% test, 75% training KDE – kernel density estimation

3) To catch up with **biology**, **2W logical neuron** can be built as **data structure** containing **model of local joint distribution**, e.g. among its connections. This way substituting some variables and normalizing would give conditional distribution, in any direction. Like in Bayes theorem but without it: straight from joint distribution. By averaging we could also use inputs as distributions, taking expected value (or median, mode) of output we can propagate values.

## **Joint distribution neurons as logical for BNNs:** multidirectional propagation, of values and probability distribution



just substitute and normalize to get conditional from joint distribution

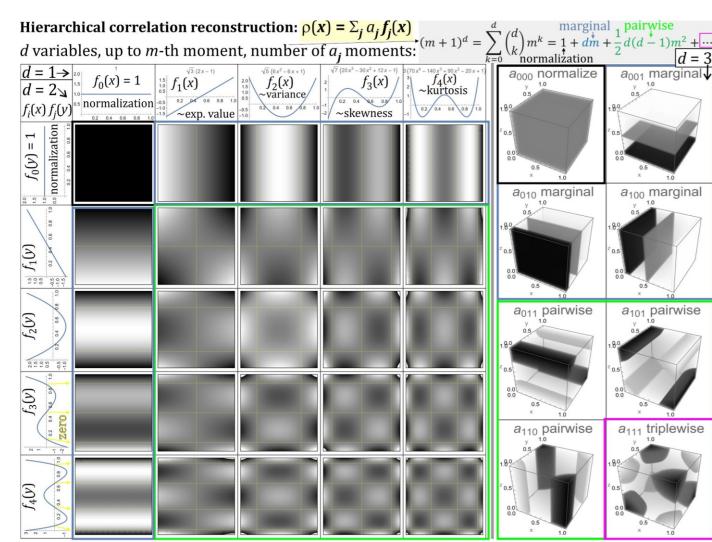
6) Basic formulas for **2D case** assuming **product basis of orthonormal** (f'(x))e.g. polynomials - with (a<sub>ii</sub>) as parameters contained by neuron with 2 connections, generally as many indexes as connections. Substituting and normalizing we get formulas for propagation as conditional distributions, or their expected values, requiring to just transpose (permutate indexes) to change propagation direction - seems simple enough to fit in biological. Working on distributions we can also control **mutual information** *I*(*X*;*Y*).

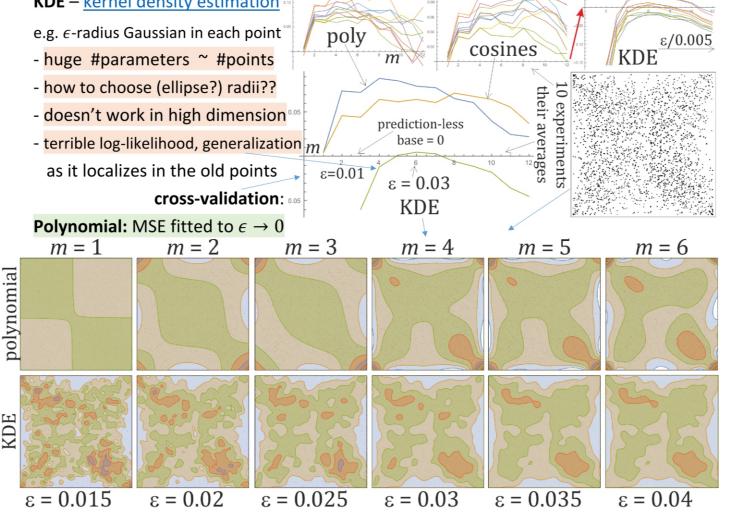
**Neuron containing local joint distribution model** – of its connections:

HIERARCHICAL CORRELATION RECONSTRUCTION  $\rho(Y,Z|X)$ for time series, conditional distribution (Bayes) models<sub>v</sub> ρ(*X*,*Y*,*Z*) (nonlinear, adaptive, all-directional) artificial neurons How to model/estimate density from a data sample? **MSE fit polynomial**  $\rho(x) = \sum_{f \in B} a_f f(x)$  (in (*f*) orthonormal basis) also for joint distribution, non-stationarity, missing data

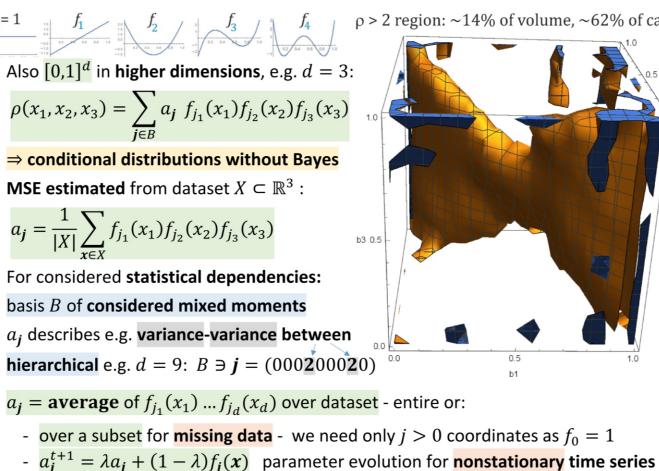
		Moments/cumulants	$\rho(x) = \sum_f a_f f(x)$	Machine learning, NN	
	# parameters	<mark>low – rough</mark>	from low to high	high - accurate	
	estimation	e.g. $m_k = \frac{1}{ X } \sum_{x \in X} x^k$	$a_f = \frac{1}{ X } \sum_{x \in X} f(x)$	usually iteration	
	Interpretable?	yes	Yes: mixed moments	depends, trust???	
	Independently?	yes	Yes (adapt, missing)	depends	
	Unique?	yes	yes (MSE)	often huge freedom	
	Accuracy?	controllable	controllable	usually uncontrollable	
	Density?	<u>moment problem</u>	<b>YES:</b> $\sum_f a_f f(x)$	depends	
	$\rightarrow$ complete	depends	yes	depends	
each variable independent ~ correlation coef.					
pair-wise $\approx$ + a11 · + a12 · + a21 · + a22 ·					

7) First [0,1] (Legendre) orthonormal polynomials and their product basis for 2,3 variables, allowing for hierarchical reconstruction of dependencies. As  $f_0=1$ , coordinates being 0 control normalization, there is no dependence with such variables. Coefficients with single nonzero index describe marginal distributions of variables, with two nonzero define pairwise dependencies, with 3 triplewise, and so on hierarchically adding higher order dependencies.





8) Analogously in *d* dimensions we need *d* index tensor representing joint distribution. Given a<sub>i</sub> coefficient describes dependencies between variables of nonzero indexes of *j* - in practice requiring to **restrict** e.g. to only pairwise dependencies up to e.g. 4th moment (kurtosis). MSE estimation by averaging can use **exponential moving average** instead for **adaptivity**. The 0 index values can be **missing**, allowing to **extract all available dependencies**.

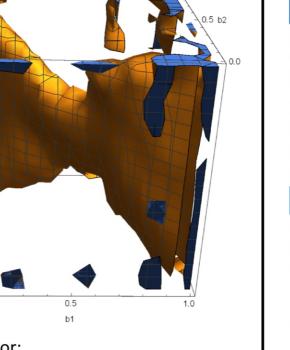


+ a<sub>11</sub>

pair-wise

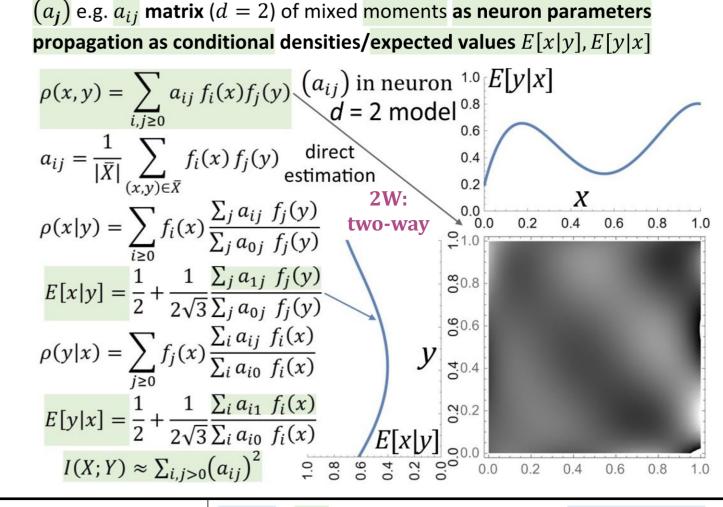
joint density

 $\rho$  > 2 region: ~14% of volume, ~62% of cases:



+ a<sub>22</sub>

var-var



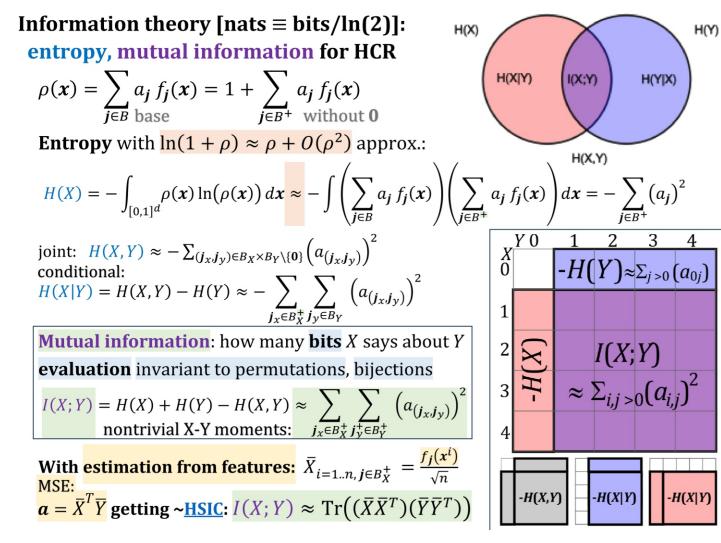
9) Can **degenerate** to very popular **KAN: Kolmogorov-Arnold networks** if only pairwise dependecies, offering many improvements.t

HCRNN: ~KAN-like parametrization, but with many advantages can propagate in any direction like biological neural networks, propagate values or probability distributions, interpretation of parameters as mixed moments, consciously add triplewise and higher dependencies, inexpensive evaluation of modeled mutual information, additional training approaches, e.g. direct estimation, tensor decomposition, information bottleneck of hidden layer

06

 $f[x_] := Exp[x[1] - x[2]^2 - x[3]^3 + x[4]^4]; d = 4; (* function, variables *$ m = 4; n = 1000; SeedRandom[1]; (\* max. features, points, next find basis: \*) fs = Expand[Table[LegendreP[i, x] \* Sqrt[(2i + 1)] /.  $x \rightarrow 2x - 1$ , {i, m}]]; X = RandomVariate[NormalDistribution[0, 1], {n, d}]; Y = Map[f, X]; inX = Table[normEDF[X[All, i]], {i, d}]; inY = normEDF[Y]; (\*normalization\*) nX = Transpose[Table[inX[i]](X[All, i]], {i, d}]; nY = inY[Y]; 4+1  $md = 4; \overline{X} = Flatten[fs[1; md]] / . x \rightarrow nX, \{\{2\}, \{3, 1\}\}]; (* features *)$ [0,1] $mu = 1; \overline{Y} = Transpose[fs[[1;;mu]] / . x \rightarrow nY]; (* only expected value *)$ as = Flatten [Transpose  $[\overline{X}] \cdot \overline{Y} ] / n;$  (\* direct parameter estimation cp = 1 / 2 + Expand [Partition [as, md].fs [[1;; md]]] / Sqrt [12]; (\* E \*) Print[Row[{ListPlot[Transpose[{nY, Y}]], "of sum of below: "}]]; joint Row[Table[cf = cp[[i]] /. x → inX[[i]][x]; (\* polynomials(normalized X)/ distr. Plot[cf, {x, -2, 2}, ImageSize  $\rightarrow$  130], {i, d}]] model KAN-like example: 1-paremeter function, summation ~exp of sum of below:  $E[x|y,z] \approx \text{normalized } \sum_{i} a_{1i0} f_i(y) + \sum_{i} a_{10i} f_i(z)$ f[x] reconstructed from n=1000 points: (normalized) 0.00.20.40.60.81.0 0.65 ~*x* 0.55 0.50 0.45 0.40 0.60 0.55 0.7

10) Neurons containing joint distribution also allow for online control of often the most valuable e.g. for training and explainability: (conditional) entropy, mutual information evaluation - in just bits of information one variable says about another, invariant to permutations and bijections. For HCR it can be **approximated** with just **sums of squares** of our moments/neuron parameters: nontrivial between variables of interest. Together with MSE estimation of (*a*) coefficients from dataset, we get simple HSIC-like formula.



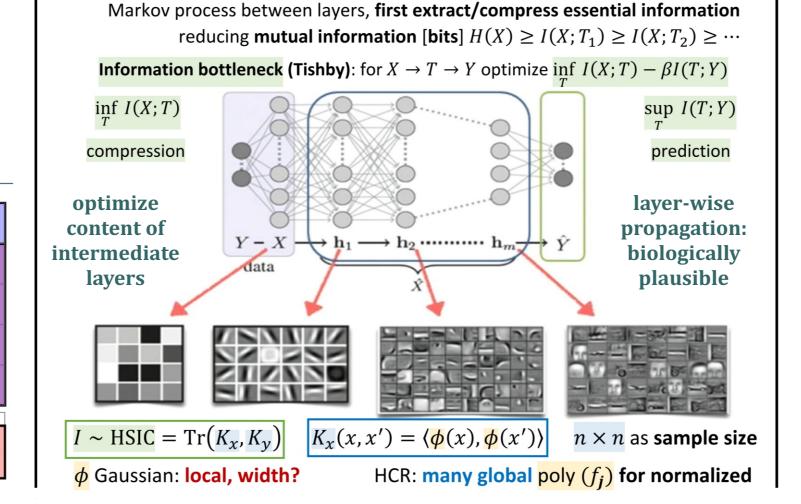
11) Biological neurons should use only local connections also for training. To improve prediction, it should maximize information about the desired **output**, but also **minimize about the input to remove noise** - allowing to optimize content based only on neighboring layers (local), preferably using mutual information as evaluation - information bottleneck approach, with formulas here as in **HSIC** (Hilbert-Schmidt Information Criterion) articles, but replacing their local basis of Gaussians, with global of polynomials/moments.

Naftali Tishby, information theoretic view - permutation, bijection independent

independent ~ correlation coef. **further statistical dependencies** 

+ a<sub>21</sub>

+ a<sub>12</sub> ·



0.30 0.2 2-2 -1 2-2 -1 -1 2-2 -2 -1

0.4

12) **Summary** with formulas for 3 variables. While single neuron can be directly trained or updated, for intermediate layers it is more difficult, but we can always treat such network as just KAN-like parametrization and train with gradient descent. One of additional training options is estimating high order tensor as dependencies, and try to automatically decompose it into smaller tensors - allowing to shift all nonlinearities to first/last layer.

0	$f_0(x) = 1 \qquad f_1(x)$	) polyn, variar	ace $f_2(x)$ ~skewness $f_4(x)$ ~kurtosis		
<b>normalization</b> ~exp. value $f_3(x)$					
NN <u>arXiv:2405.05097</u>	d=3 variables	basis in [0,1]	$f_0 = 1, f_1 \propto 2x - 1, \ \int_0^1 f_i(x) f_j(x) dx = \delta_{ij}$		
Reduced to ~KAN for	<b>HCR</b> neuron	↓↑ normalize	$x \leftrightarrow \text{CDF}(x) \sim U[0,1]$ empirical/param.		
pairwise-only dependencies	$\begin{array}{c} & \text{normalize } CDF \\ & -1 + calculate \\ CDF & \{f\} \text{ or } \{g\} \\ & g_{ijk} \\ \end{array}$	HCR joint density static estimation	$\rho(x, y, z) = \sum_{ijk \in B} a_{ijk} f_i(x) f_j(y) f_k(z)$		
Additionally:		from $\overline{X}$ dataset	mean: $a_{ijk} = \frac{1}{ \overline{X} } \sum_{(x,y,z) \in \overline{X}} f_i(x) f_j(y) f_k(z)$		
		dynamic (EMA) model <b>update</b>	$a_{ijk} \xrightarrow{(x,y,z)} (1-\lambda)a_{ijk} + \lambda f_i(x)f_j(y)f_k(z)$		
can extend to $\gamma$	ρ	$\rho(X = x   y, z) \approx$			
triplewise or higher,		1 conditional	$\sum_{i} f_i(x) \frac{\sum_{jk} a_{ijk} f_j(y) f_k(z)}{\sum_{jk} a_{0jk} f_j(y) f_k(z)} \qquad \frac{\text{current}}{\text{normal.}}$		
omnidirectional		$E[X = x y, z] \approx$	$\frac{1}{2} + \frac{1}{2\sqrt{3}} \frac{\sum_{jk} a_{1jk} f_j(y) f_k(z)}{\sum_{jk} a_{0jk} f_j(y) f_k(z)} \qquad \text{sufficient} \\ \text{if norm.}$		
propagation $\uparrow \leftrightarrow$ ,		↑ propagation?			
direct $a_j$ parameter	Z	$ \rho(y, z   x) \approx $ $\downarrow \text{ conditional}$	$\sum_{ik} f_j(y) f_k(z) \frac{\sum_i a_{ijk} f_i(x)}{\sum_i a_{i00} f_i(x)} \qquad \frac{\text{current}}{\text{normal.}}$		
estimation/update,	$\rho(\mathbf{x}) = \sum_{j} a_{j} f_{j}(\mathbf{x})$ $\rho(\mathbf{x}, \mathbf{y}, \mathbf{z}) \text{ for } d = 3$	$E[Y = y x] \approx$ \$\propagation?	$\frac{1}{2} + \frac{1}{2\sqrt{3}} \frac{\sum_{j} a_{1j0} f_{j}(y)}{\sum_{j} a_{0j0} f_{j}(y)} \stackrel{\text{polyn.}}{\underset{\text{like if normalized}}{\text{polyn.}} \text{sufficient}$		
propagate values or	$\rho(x y,z) \text{ for } a = 3$ propagation: $\rho(x y,z)$ $\rho(y,z x)$	entropy,	$H(X) \approx -\sum_{j \in B_{\mathbf{v}}^+} (a_j)^2  \text{[nits]}$		
probability distributions,		mutual information	$I(X;Y) \approx \sum_{j_X \in B_X^+} \sum_{j_Y \in B_Y^+} \left(a_{(j_X,j_Y)}\right)^2$		
interpretation: moments,	pairwise	basis optimization	$M_{i,jk} = a_{ijk}$ SVD: $MM^T = \sum_i \sigma_i  \boldsymbol{v}_i \boldsymbol{v}_i^T$		
- cheaply calculate entropy	$\sim \mathbf{KAN}: i \cdot j \cdot k = 0$	$(y,z) \rightarrow x$ $\{f_i(x)\} \rightarrow \{g_i(x)\}$	$g_i(x) = \sum_j v_{ij} f_j(x) \qquad v_i \cdot v_j = \delta_{ij}$ $f_i = \sum_l v_{li} g_l \qquad a_{ijk} \to \sum_l v_{li} a_{ljk}$		
mutual information,	intermediate layers/variables?				
additional training e.g.	<b>neural</b> <b>network</b> - standard <b>backpropagation</b> of $a_{ijk}$ gradients - <b>Information bottleneck method</b> for neurons				
tensor decomposition,	$u \operatorname{or} \int \rho(u) = j$ k = -up/down propagation				
information bottleneck <i>a b + a<sub>ijk</sub></i> estimation/update					
$B_{\rm Y}^+$ - tensor decomposition					
$B_X^+ \bullet \bullet \bullet T?$			$A_{i_1 i_2 i_3 i_4}^k \approx \sum_{j_1, j_2} a_{i_1 i_2}^{j_1} b_{i_3 i_4}^{j_2} c_{j_1 j_2}^k$		