

Neural Encoding Models

Maneesh Sahani

**Gatsby Computational Neuroscience Unit
University College London**

November 2015

Neural Coding

The brain appears to be modular. Different structures and cortical areas compute, represent and transmit separate pieces of information.

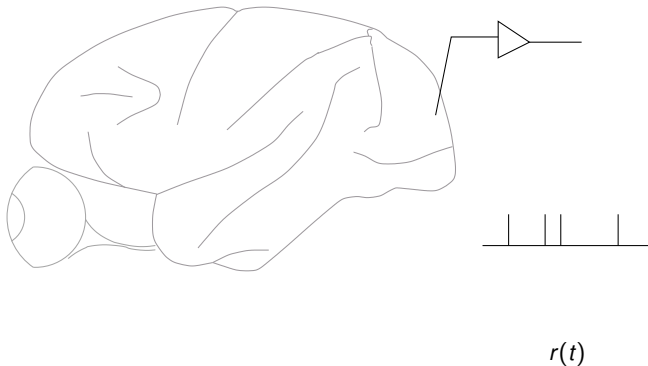
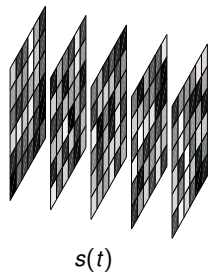
The coding questions:

- ▶ What information is represented by a particular neural population?
 - ▶ easy (?) if we know the code
 - ▶ more generally, can search for selectivity / invariance
 - ▶ encoded quantities might not be obvious: inferred latent variables, uncertainty . . .
- ▶ How is that information encoded?
 - ▶ firing rate, spiking timing (relative to other spikes, population oscillations, onset of time-invariant stimulus)?
 - ▶ functional mapping of encoded variable to spikes?
 - ▶ easy (?) if we know what is encoded

A complete answer will require convergence of theory and empirical results.

Computation plays a vital part in systematising empirical data.

Stimulus coding



Decoding: $\hat{s}(t) = G[r(t)]$

(reconstruction)

Encoding: $\hat{r}(t) = F[s(t)]$

(systems identification)

Why?

The stimulus coding problem has sometimes been identified with the “neural coding” problem.

However, on the face of it, mapping *either* the decoding or encoding function does not by itself answer either of our basic questions about coding.

So why do we do it?

- ▶ encapsulate and systematise the response so that we *can* ask the questions that we want answered.
- ▶ design hypothesis-driven stimulus-coding models: evaluate coding reliability for different function(al)s of $s(t)$ and for different definitions of $r(t)$.
- ▶ but correlation \nrightarrow causation: in this case the *presence* of information about an aspect of the stimulus in a particular aspect of the response does not mean that the brain *uses* that information.

General approach

Goal: Estimate $p(\text{spike}|s, H)$ [or $\lambda(t|s[0, t], H(t))$] from data.

- ▶ Naive approach: measure $p(\text{spike}, H|s)$ directly for every setting of s .
 - ▶ too hard: too little data and too many potential inputs.
- ▶ Estimate some functional $F[\rho]$ instead (e.g. mutual information)
- ▶ Select stimuli efficiently
- ▶ **Fit models with smaller numbers of parameters**

Spikes, or rate?

Most neurons communicate using action potentials — statistically described by a **point process**:

$$P(\text{spike} \in [t, t + dt]) = \lambda(t|H(t), \text{stimulus}, \text{network activity})dt$$

To fully model the response we need to identify λ . In general this depends on spike history $H(t)$ and network activity. Three options:

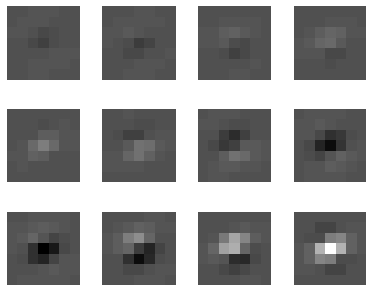
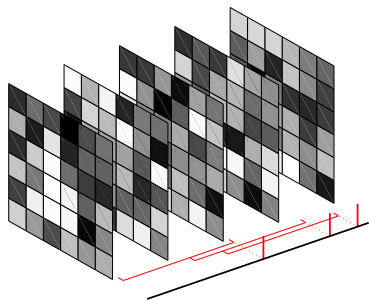
- ▶ Ignore the history dependence, take network activity as source of “noise” (i.e. assume firing is inhomogeneous Poisson or Cox process, conditioned on the stimulus).
- ▶ Average multiple trials to estimate the mean intensity (or PSTH)

$$\bar{\lambda}(t, \text{stimulus}) = \lim_{N \rightarrow \infty} \frac{1}{N} \sum_n \lambda(t|H_n(t), \text{stimulus}, \text{network}_n),$$

and try to fit this.

- ▶ Attempt to capture history and network effects in simple models.

Spike-triggered average



Decoding:

mean of $P(s | r = 1)$

Encoding:

predictive filter

Linear regression

$$r(t) = \int_0^T s(t-\tau)w(\tau)d\tau$$

s_1	s_2	s_3	\dots	s_T	s_{T+1}	\dots
-------	-------	-------	---------	-------	-----------	---------

$\underbrace{\hspace{10em}}$
 $\underbrace{\hspace{10em}}$

s_1	s_2	s_3	\dots	s_T
s_2	s_3	s_4	\dots	s_{T+1}
		\vdots		

 \times

w_t
\vdots
w_3
w_2
w_1

 $=$

r_T
r_{T+1}
\vdots

$$SW = R$$

$$W(\omega) = \frac{S(\omega)^* R(\omega)}{|S(\omega)|^2}$$

$$W = \underbrace{(S^T S)^{-1}}_{\Sigma_{SS}} \underbrace{(S^T R)}_{STA}$$

Linear models

So the (whitened) spike-triggered average gives the minimum-squared-error linear model.

Issues:

- ▶ overfitting and regularisation
 - ▶ standard methods for regression
- ▶ negative predicted rates
 - ▶ can model deviations from background
- ▶ real neurons aren't linear
 - ▶ models are still used extensively
 - ▶ interpretable suggestions of underlying sensitivity (but see later)
 - ▶ may provide unbiased estimates of cascade filters (see later)

How good are linear predictions?

We would like an absolute measure of model performance. Two things make this difficult:

Measured responses can never be predicted perfectly, even in principle:

- ▶ The measurements themselves are noisy.

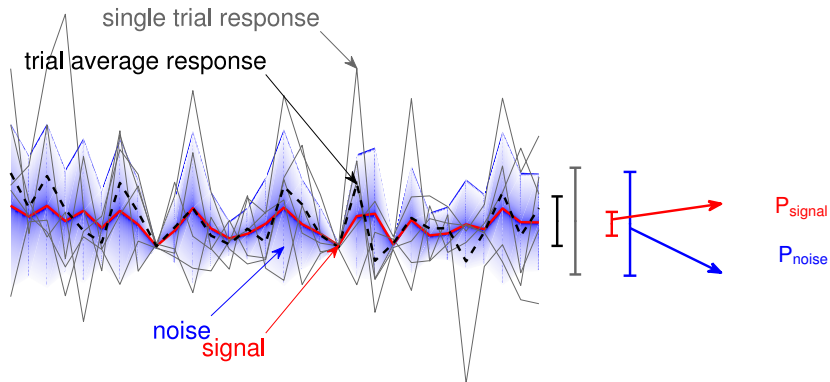
Even if we can discount this, a model may predict poorly because either:

- ▶ It is the wrong model.
- ▶ The parameters are mis-estimated due to noise.

Approaches:

- ▶ Compare $I(\text{resp}; \text{pred})$ to $I(\text{resp}; \text{stim})$.
 - ▶ mutual information estimators are biased
- ▶ Compare $E(\text{resp} - \text{pred})$ to $E(\text{resp} - \text{psth})$ where psth is gathered over a very large number of trials.
 - ▶ may require impractical amounts of data to estimate the psth
- ▶ Compare the *predictive power* to the *predicable power* (similar to ANOVA).

Estimating predictable power



$$\underbrace{\text{response}}_{\mathbf{r}^{(n)}} = \text{signal} + \text{noise}$$

$$\left. \begin{aligned} \overline{P(\mathbf{r}^{(n)})} &= P_{\text{signal}} + P_{\text{noise}} \\ P(\overline{\mathbf{r}^{(n)}}) &= P_{\text{signal}} + \frac{1}{N} P_{\text{noise}} \end{aligned} \right\} \Rightarrow \begin{cases} \hat{P}_{\text{signal}} = \frac{1}{N-1} \left(NP(\overline{\mathbf{r}^{(n)}}) - \overline{P(\mathbf{r}^{(n)})} \right) \\ \hat{P}_{\text{noise}} = \overline{P(\mathbf{r}^{(n)})} - \hat{P}_{\text{signal}} \end{cases}$$

Testing a model

For a perfect prediction

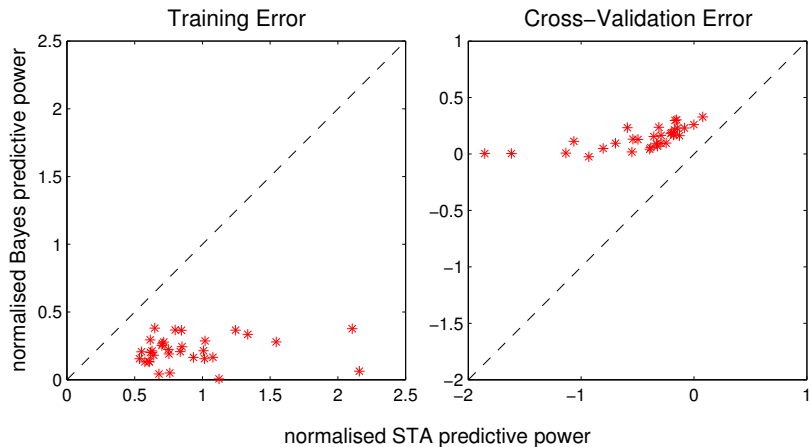
$$\langle P(\overline{\text{trial}}) - P(\text{residual}) \rangle = P(\text{signal})$$

Thus, we can judge the performance of a model by the **normalized predictive power**

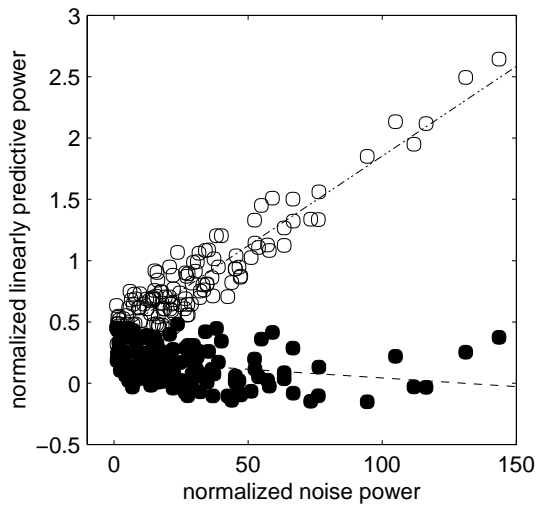
$$\frac{P(\overline{\text{trial}}) - P(\text{residual})}{\hat{P}(\text{signal})}$$

Similar to coefficient of determination (r^2), but the denominator is the **predictable** variance.

Predictive performance



Extrapolating the model performance



Jackknife bias correction

Estimate bias by extrapolation in data size:

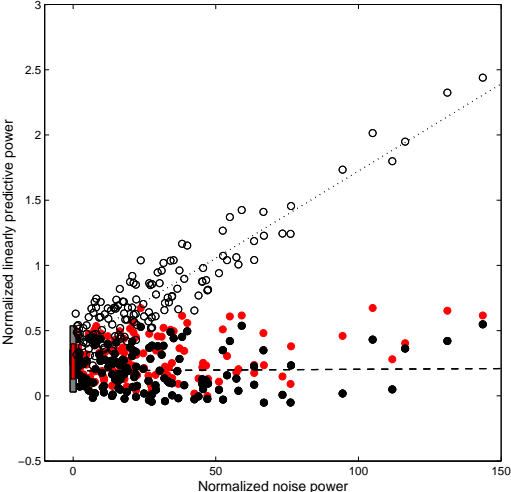
$$\mathcal{T}_{\text{jn}} = N\mathcal{T} - (N - 1)\mathcal{T}_{100}$$

where \mathcal{T} is the training error on all data and \mathcal{T}_{100} is the average training error on all sets of $N - 1$ data.

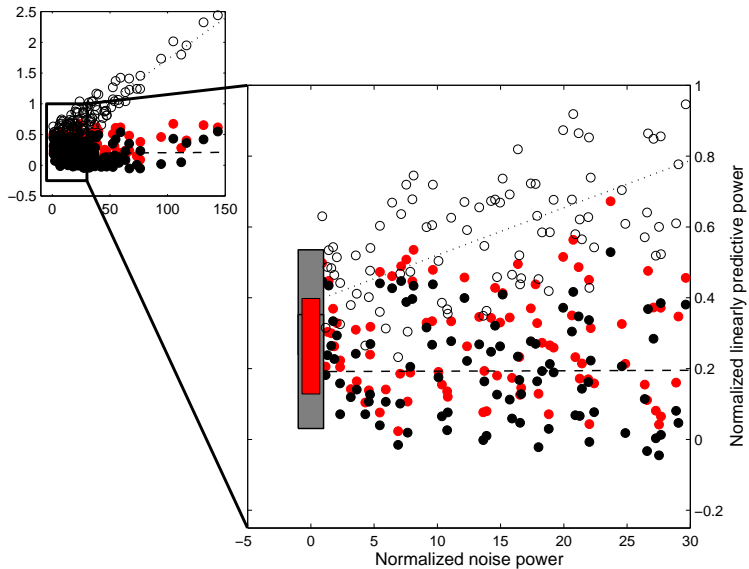
For a linear model we can find this in closed form:

$$\mathcal{T}_{\text{jn}} = \frac{1}{N} \sum_i \left(\frac{(r_i - \mathbf{s}_i \mathbf{w}^{\text{ML}})^2}{1 - \mathbf{s}_i (\mathbf{S}^T \mathbf{S})^{-1} \mathbf{s}_i^T} \right)$$

Jackknifed estimates

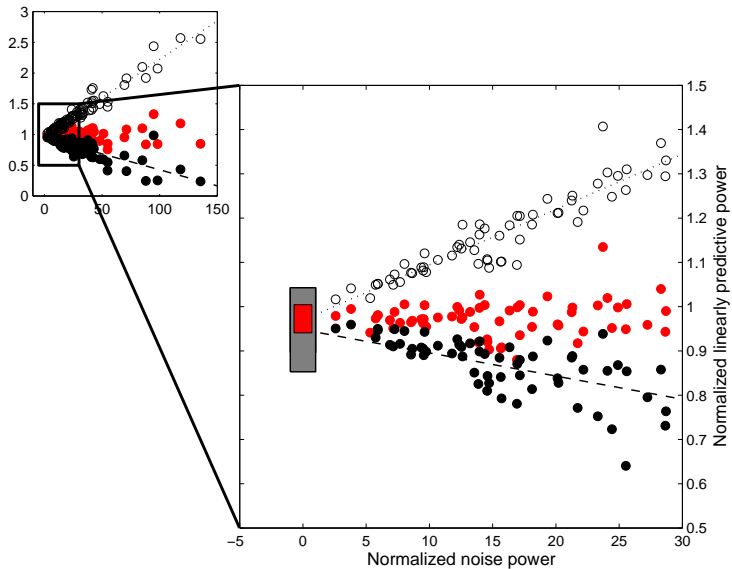


Extrapolated linearity



[extrapolated range: (0.19,0.39); mean Jackknife estimate: 0.29]

Simulated (almost) linear data



[extrapolated range: (0.95,0.97); mean Jackknife estimate: 0.97]

Beyond linearity

Beyond linearity

Linear models often fail to predict well. Alternatives?

- ▶ Wiener/Volterra functional expansions
 - ▶ M-series
 - ▶ Linearised estimation
 - ▶ Kernel formulations
- ▶ LN (Wiener) cascades
 - ▶ Spike-trigger covariance (STC) methods
 - ▶ “Maximally informative” dimensions (MID) \Leftrightarrow ML nonparametric LNP models
 - ▶ ML Parametric GLM models
- ▶ NL (Hammerstein) cascades
 - ▶ Multilinear formulations

The Volterra functional expansion

A polynomial-like expansion for functionals (or operators).

Let $y(t) = F[x(t)]$. Then:

$$y(t) \approx k^{(0)} + \int d\tau k^{(1)}(\tau)x(t-\tau) + \iint d\tau_1 d\tau_2 k^{(2)}(\tau_1, \tau_2)x(t-\tau_1)x(t-\tau_2) \\ + \iiint d\tau_1 d\tau_2 d\tau_3 k^{(3)}(\tau_1, \tau_2, \tau_3)x(t-\tau_1)x(t-\tau_2)x(t-\tau_3) + \dots$$

or (in discretised time)

$$y_t = K^{(0)} + \sum_i K_i^{(1)} x_{t-i} + \sum_{ij} K_{ij}^{(2)} x_{t-i} x_{t-j} + \sum_{ijk} K_{ijk}^{(3)} x_{t-i} x_{t-j} x_{t-k} + \dots$$

For finite expansion, the kernels $k^{(0)}$, $k^{(1)}(\cdot)$, $k^{(2)}(\cdot, \cdot)$, $k^{(3)}(\cdot, \cdot, \cdot)$, ... are not straightforwardly related to the functional F . Indeed, values of lower-order kernels change as the maximum order of the expansion is increased.

Estimation: model is linear in kernels, so can be estimated just like a linear (first-order) model with expanded “input”.

- ▶ Kernel trick: polynomial kernel $K(x_1, x_2) = (1 + x_1 x_2)^n$.
- ▶ M-series.

Wiener Expansion

The Wiener expansion gives functionals of different orders that are **orthogonal** for white noise input $x(t)$.

$$G_0[x(t); h^{(0)}] = h^{(0)}$$

$$G_1[x(t); h^{(1)}] = \int dx d\tau h^{(1)}(\tau) x(t - \tau)$$

$$G_2[x(t); h^{(2)}] = \iint d\tau_1 d\tau_2 h^{(2)}(\tau_1, \tau_2) x(t - \tau_1) x(t - \tau_2) - P \int dx d\tau_1 h^{(2)}(\tau_1, \tau_1)$$

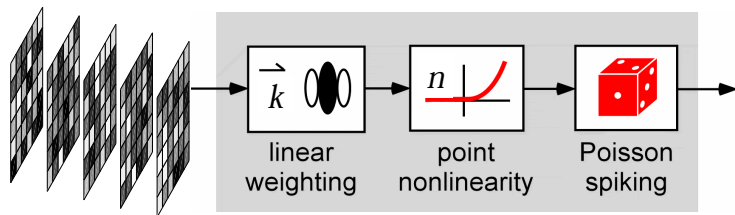
$$G_3[x(t); h^{(3)}] = \iiint d\tau_1 d\tau_2 d\tau_3 h^{(3)}(\tau_1, \tau_2, \tau_3) x(t - \tau_1) x(t - \tau_2) x(t - \tau_3) \\ - 3P \iint d\tau_1 d\tau_2 h^{(3)}(\tau_1, \tau_2, \tau_2) x(t - \tau_1)$$

Easy to verify that $\mathbb{E}[G_i[x(t)]G_j[x(t)]] = 0$ for $i \neq j$.

Thus, these kernels can be estimated independently. **But, they depend on the stimulus.**

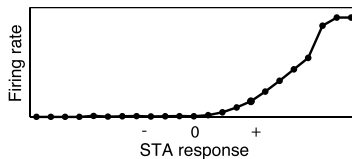
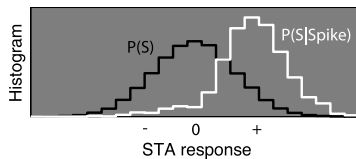
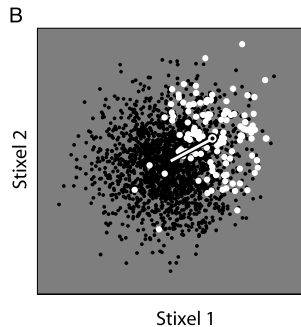
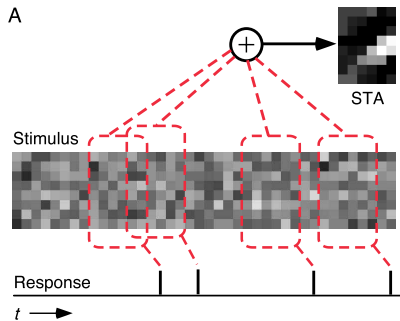
Cascade models

The LNP (Wiener) cascade

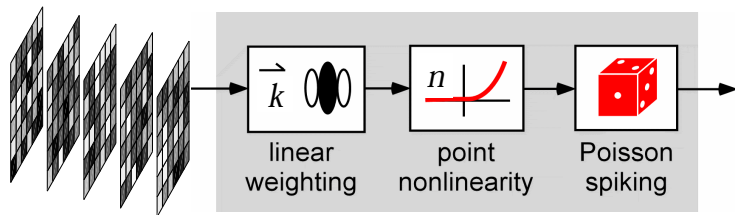


- ▶ Rectification addresses negative firing rates.
- ▶ Loose biophysical correspondence.

LNP estimation – the Spike-triggered ensemble

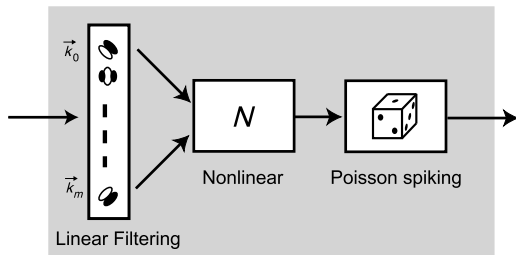


Single linear filter



- ▶ STA is **unbiased** estimate of filter for spherical input distribution. (Bussgang's theorem)
- ▶ Elliptically-distributed data can be whitened \Rightarrow linear regression weights are **unbiased**.
- ▶ Linear weights are not necessarily maximum-likelihood (or otherwise optimal), even for spherical/elliptical stimulus distributions.
- ▶ Linear weights may be biased for general stimuli (binary/uniform or natural).

Multiple filters

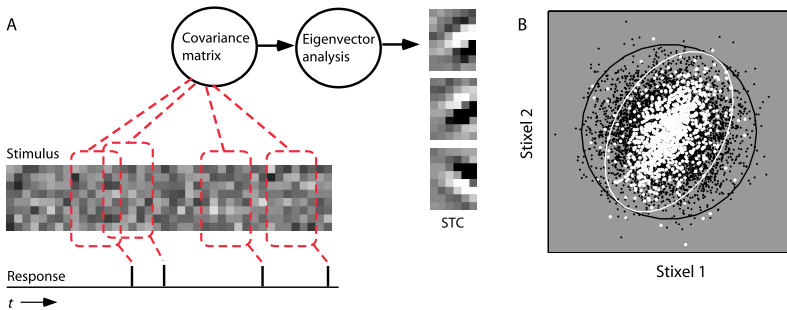


Distribution changes along relevant directions (and, usually, along all linear combinations of relevant directions).

Proxies to measure change in distribution:

- ▶ mean: STA (can only reveal a single direction)
- ▶ variance: STC
- ▶ binned (or kernel) KL divergence: MID “maximally informative directions” (equivalent to ML in LNP model with binned nonlinearity)

STC



Project out STA:

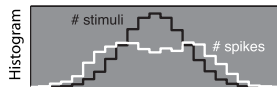
$$\tilde{X} = X - (X\mathbf{k}_{\text{sta}})\mathbf{k}_{\text{sta}}^T; \quad C_{\text{prior}} = \frac{\tilde{X}^T \tilde{X}}{N}; \quad C_{\text{spike}} = \frac{\tilde{X}^T \text{diag}(Y)\tilde{X}}{N_{\text{spike}}}$$

Choose directions with greatest change in variance:

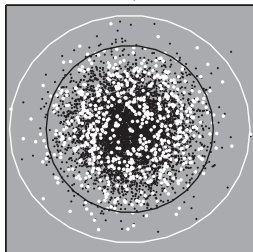
$$\underset{\|\mathbf{v}\|=1}{\text{k-argmax}} \mathbf{v}^T (C_{\text{prior}} - C_{\text{spike}}) \mathbf{v}$$

⇒ find eigenvectors of $(C_{\text{prior}} - C_{\text{spike}})$ with large (absolute) eigvals.

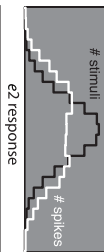
Reconstruct nonlinearity (may assume separability)



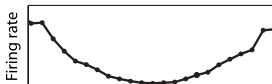
e1 response



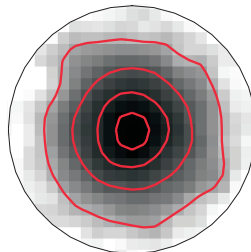
Histogram



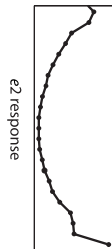
e2 response



e1 response



Firing rate



Biases

STC (obviously) requires that the nonlinearity alter variance.

If so, subspace is unbiased provided distribution is

- ▶ radially (elliptically) symmetric
- ▶ AND independent

⇒ Gaussian.

May be possible to correct for non-Gaussian stimulus by transformation, subsampling or weighting (latter two at cost of variance).

More LNP methods

- ▶ Non-parametric non-linearities:

“Maximally informative dimensions” (MID) \Leftrightarrow “non-parametric” maximum likelihood.

- ▶ Intuitively, extends the variance difference idea to arbitrary differences between marginal and spike-conditioned stimulus distributions.

$$\mathbf{k}_{\text{MID}} = \underset{\mathbf{k}}{\operatorname{argmax}} \operatorname{KL}[P(\mathbf{k} \cdot \mathbf{x}) || P(\mathbf{k} \cdot \mathbf{x} | \text{spike})]$$

- ▶ Measuring KL requires binning or smoothing—turns out to be equivalent to fitting a non-parametric nonlinearity by binning or smoothing.
 - ▶ Difficult to use for high-dimensional LNP models (but ML viewpoint suggests separable or “cylindrical” basis functions).
-
- ▶ Parametric non-linearities: the “generalised linear model” (GLM).

Generalised linear models

LN models with specified nonlinearities and exponential-family noise.

In general (for monotonic g):

$$y \sim \text{ExpFamily}[\mu(\mathbf{x})]; \quad g(\mu) = \beta\mathbf{x}$$

For our purposes easier to write

$$y \sim \text{ExpFamily}[f(\beta\mathbf{x})]$$

(Continuous time) point process likelihood with GLM-like dependence of λ on covariates is approached in limit of bins $\rightarrow 0$ by either Poisson or Bernoulli GLM.

Mark Berman and T. Rolf Turner (1992) Approximating Point Process Likelihoods with GLIM
Journal of the Royal Statistical Society. Series C (Applied Statistics), 41(1):31-38.

Generalised linear models

Poisson distribution $\Rightarrow f = \exp()$ is *canonical* (*natural params* = $\beta\mathbf{x}$).

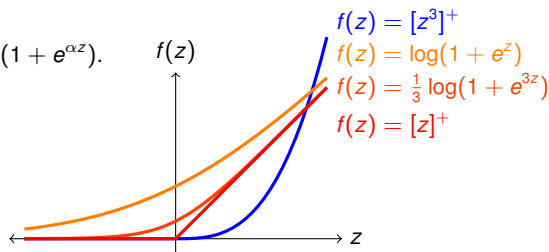
Canonical link functions give concave likelihoods \Rightarrow unique maxima.

Generalises (for Poisson) to any f which is convex and log-concave:

$$\text{log-likelihood} = c - f(\beta\mathbf{x}) + y \log f(\beta\mathbf{x})$$

Includes:

- ▶ threshold-linear
- ▶ threshold-polynomial
- ▶ “soft-threshold” $f(z) = \alpha^{-1} \log(1 + e^{\alpha z})$.



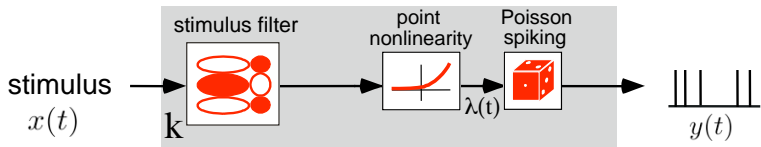
Generalised linear models

ML parameters found by

- ▶ gradient ascent
- ▶ IRLS

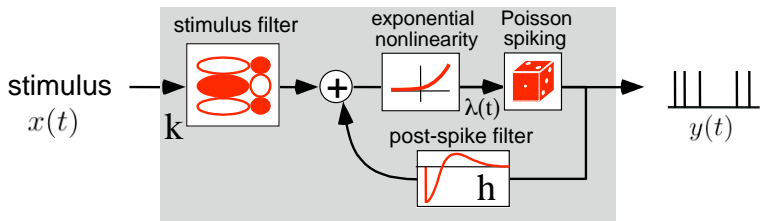
Regularisation by L_2 (quadratic) or L_1 (absolute value – sparse) penalties (MAP with Gaussian/Laplacian priors) preserves concavity.

Linear-Nonlinear-Poisson (GLM)



GLM with history-dependence

(Truccolo et al 04)

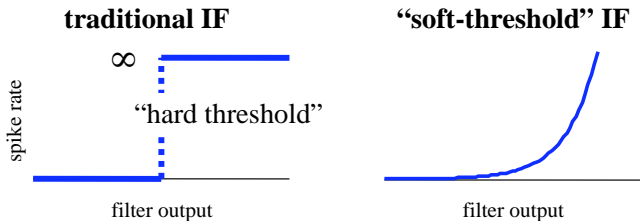
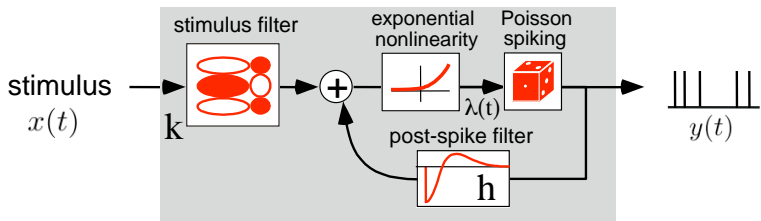


conditional intensity (spike rate)

$$\lambda(t) = f(k \cdot x(t) + h \cdot y(t))$$
$$= e^{k \cdot x(t)} \cdot e^{h \cdot y(t)}$$

- rate is a product of stim- and spike-history dependent terms
- output no longer a Poisson process
- also known as “soft-threshold” Integrate-and-Fire model

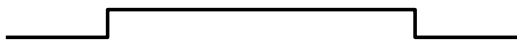
GLM with history-dependence



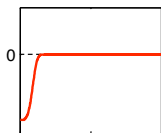
- "soft-threshold" approximation to Integrate-and-Fire model

GLM dynamic behaviors

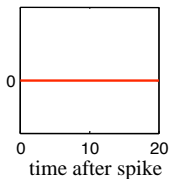
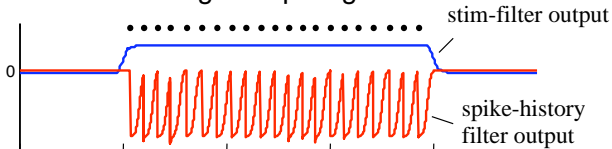
stimulus $x(t)$



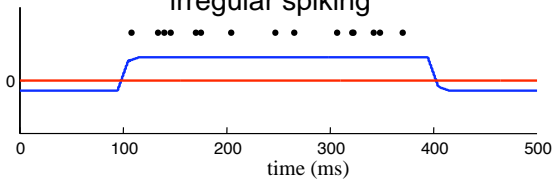
post-spike waveform



regular spiking

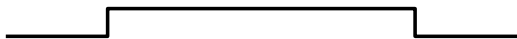


irregular spiking

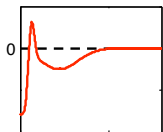


GLM dynamic behaviors

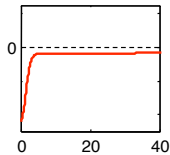
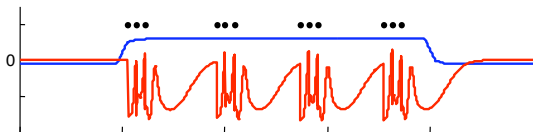
stimulus $x(t)$



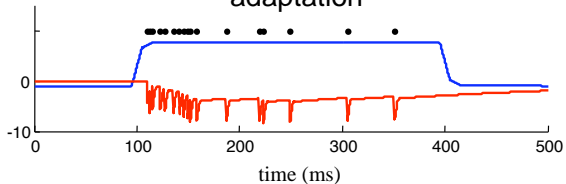
post-spike waveform



bursting



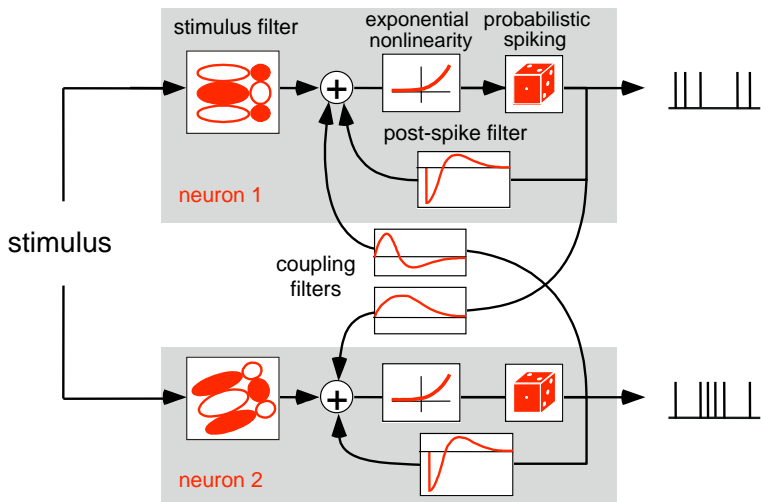
adaptation



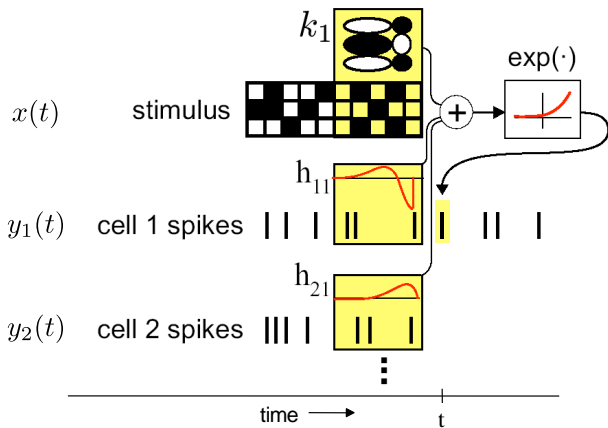
time after spike

time (ms)

multi-neuron GLM



GLM equivalent diagram:



conditional intensity
(spike rate)

$$\lambda_i(t) = \exp(k_i \cdot x(t) + \sum_j h_{ij} \cdot y_j(t))$$

Non-LN models?

The idea of responses depending on one or a few linear stimulus projections has been dominant, but cannot capture all non-linearities.

- ▶ Contrast sensitivity might require normalisation by $\|\mathbf{s}\|$.
- ▶ Linear weighting may depend on *units* of stimulus measurement: amplitude? energy? logarithms? thresholds? (NL models – Hammerstein cascades)
- ▶ Neurons, particularly in the auditory system are known to be sensitive to combinations of inputs: forward suppression; spectral patterns (Young); time-frequency interactions (Sadogopan and Wang).
- ▶ Experiments with realistic stimuli reveal nonlinear sensitivity to parts/whole (Bar-Yosef and Nelken).

Many of these questions can be tackled using a multilinear (cartesian tensor) framework.

Input nonlinearities

The basic linear model (for sounds):

$$\underbrace{\hat{r}(i)}_{\text{predicted rate}} = \sum_{jk} \underbrace{w_{jk}^{\text{tf}}}_{\text{STRF weights}} \underbrace{s(i-j, k)}_{\text{stimulus power}},$$

How to measure s ? (pressure, intensity, dB, thresholded, ...)

We can *learn* an optimal representation $g(\cdot)$:

$$\hat{r}(i) = \sum_{jk} w_{jk}^{\text{tf}} g(s(i-j, k)).$$

Define: basis functions $\{g_l\}$ such that $g(s) = \sum_l w_l^l g_l(s)$
and stimulus array $M_{ijkl} = g_l(s(i-j, k))$. Now the model is

$$\hat{r}(i) = \sum_j w_{jk}^{\text{tf}} w_l^l M_{ijkl} \quad \text{or} \quad \hat{\mathbf{r}} = (\mathbf{w}^{\text{tf}} \otimes \mathbf{w}^l) \bullet \mathbf{M}.$$

Multilinear models

Multilinear forms are straightforward to optimise by alternating least squares.

Cost function:

$$\mathcal{E} = \left\| \mathbf{r} - (\mathbf{w}^{\text{tf}} \otimes \mathbf{w}^{\text{l}}) \bullet \mathbf{M} \right\|^2$$

Minimise iteratively, defining *matrices*

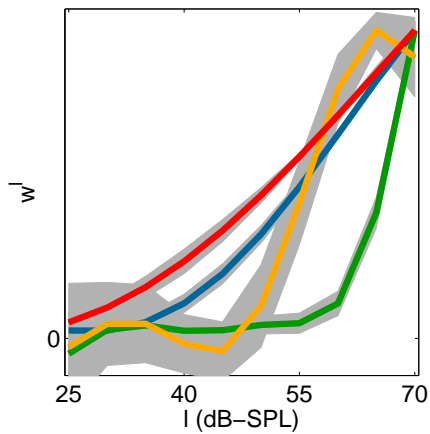
$$\mathbf{B} = \mathbf{w}^{\text{l}} \bullet \mathbf{M} \quad \text{and} \quad \mathbf{A} = \mathbf{w}^{\text{tf}} \bullet \mathbf{M}$$

and updating

$$\mathbf{w}^{\text{tf}} = (\mathbf{B}^{\text{T}} \mathbf{B})^{-1} \mathbf{B}^{\text{T}} \mathbf{r} \quad \text{and} \quad \mathbf{w}^{\text{l}} = (\mathbf{A}^{\text{T}} \mathbf{A})^{-1} \mathbf{A}^{\text{T}} \mathbf{r}.$$

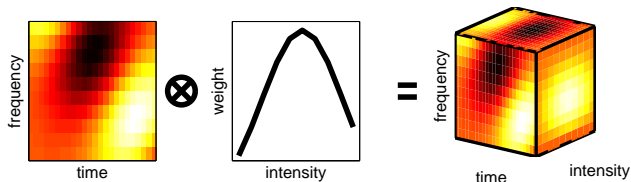
Each linear regression step can be regularised by evidence optimisation (suboptimal), with uncertainty propagated approximately using *variational* methods.

Some input non-linearities



Parameter grouping

Separable models: (time) \otimes (frequency). The input nonlinearity model is separable in another sense: (time, frequency) \otimes (sound level).

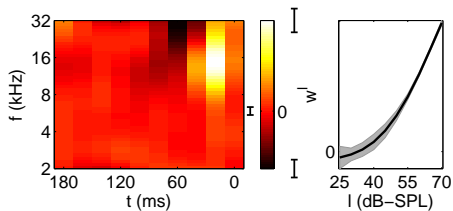


Other separations:

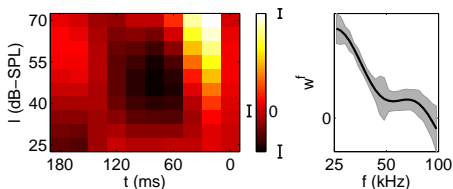
- ▶ (time, sound level) \otimes (frequency): $\hat{\mathbf{r}} = (\mathbf{w}^{\text{tl}} \otimes \mathbf{w}^f) \bullet \mathbf{M}$,
- ▶ (frequency, sound level) \otimes (time): $\hat{\mathbf{r}} = (\mathbf{w}^{\text{fl}} \otimes \mathbf{w}^t) \bullet \mathbf{M}$,
- ▶ (time) \otimes (frequency) \otimes (sound level): $\hat{\mathbf{r}} = (\mathbf{w}^l \otimes \mathbf{w}^f \otimes \mathbf{w}^l) \bullet \mathbf{M}$.

Some examples

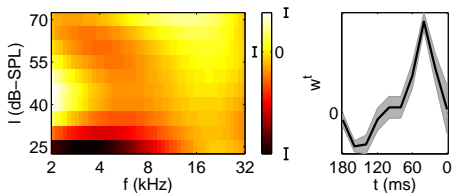
(time, frequency) \otimes (sound level):



(time, sound level) \otimes (frequency):



(frequency, sound level) \otimes (time):

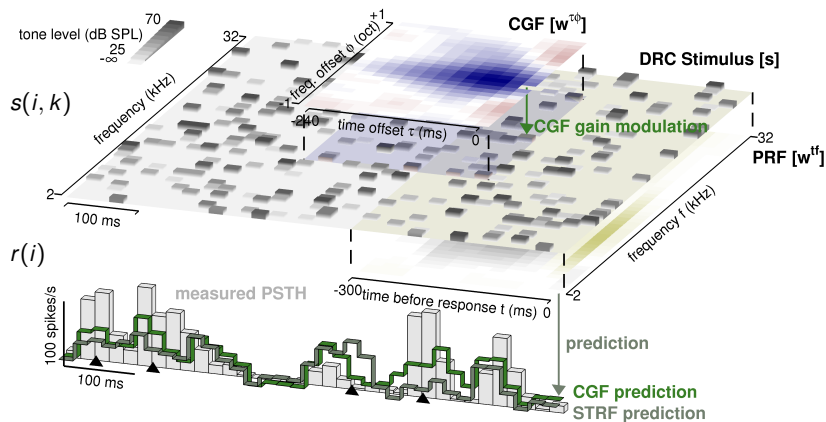


Variable (combination-dependent) input gain

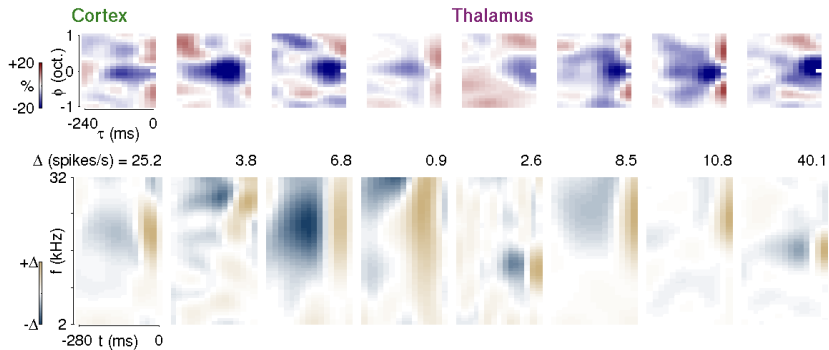
- ▶ Sensitivities to different points in sensory space are not independent.
- ▶ Rather, the sensitivity at one point depends on other elements of the stimulus that create a *local* sensory context.
- ▶ This context adjusts the **input gain** of the cell from moment to moment, dynamically refining the shape of the weighted receptive field.

A context-sensitive model

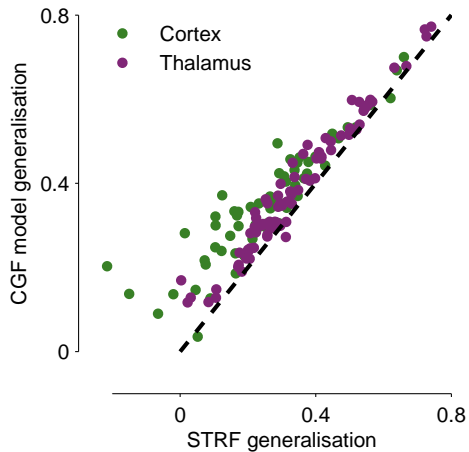
$$\hat{r}(i) = c + \sum_{j=0}^J \sum_{k=1}^K w_{j+1,k}^{f\phi} s(i-j, k) \left(1 + \sum_{m=0}^M \sum_{n=-N}^N w_{m+1,n+N+1}^{\tau\phi} s(i-j-m, k+n) \right)$$



Some examples

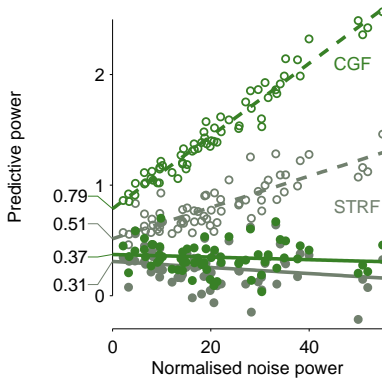


Predictive performance

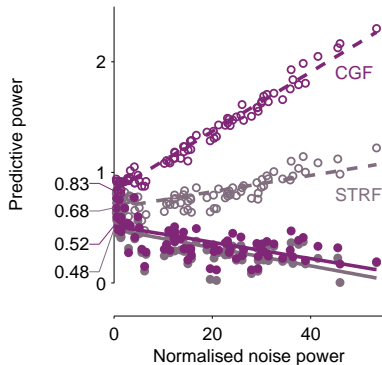


Predictive performance

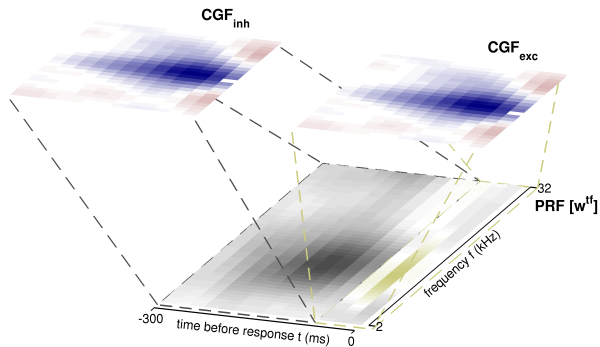
Cortex



Thalamus



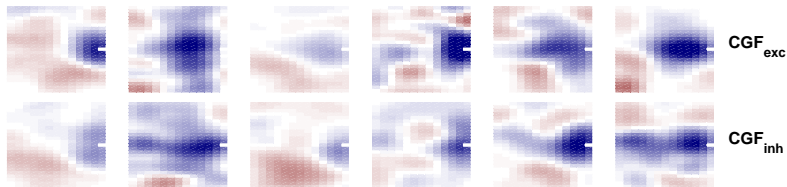
CGF consistency across the PRF



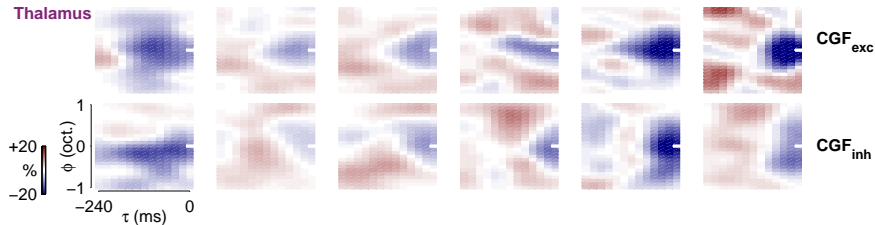
- ▶ As the CGF can be associated with the PRF weights rather than the stimulus, we can apply different CGFs to different PRF domains.

CGF consistency across the PRF

Cortex

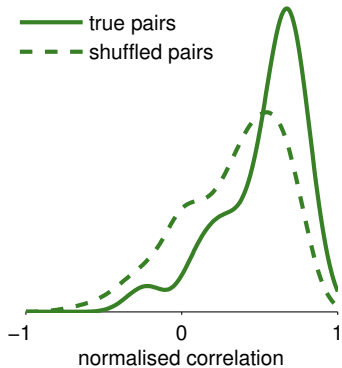


Thalamus

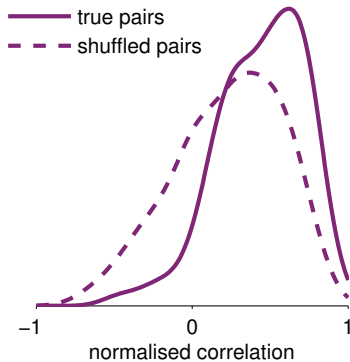


CGF consistency across the PRF

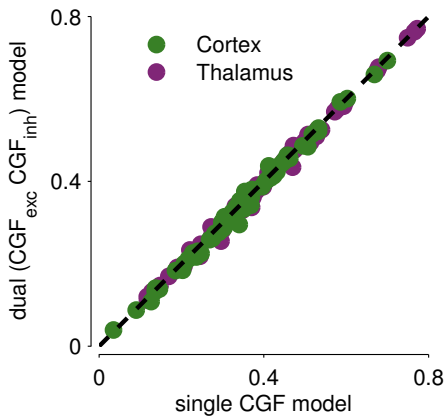
Cortex



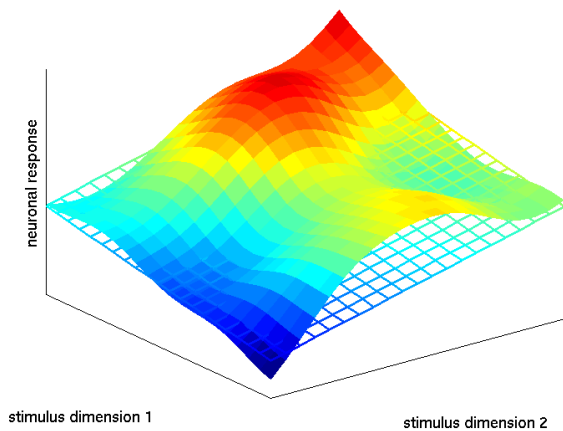
Thalamus



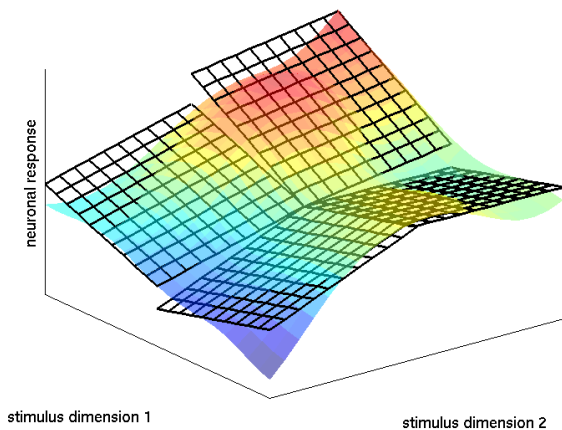
CGF consistency across the PRF



Linear fits to non-linear functions



Approximations are stimulus dependent



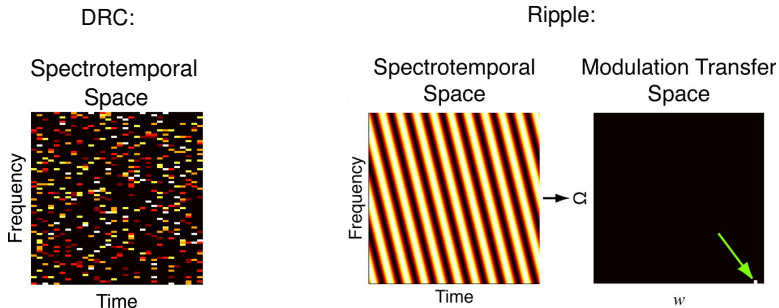
(Stimulus dependence does not always signal response adaptation)

Consequences

Local fitting can have counterintuitive consequences on the interpretation of a “receptive field”.

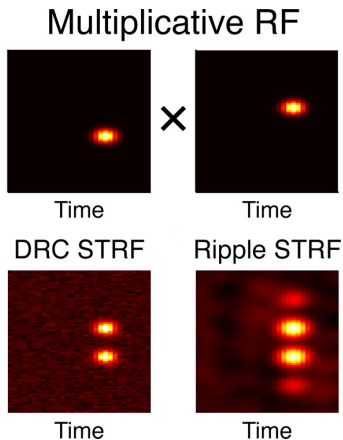
“Independently distributed” stimuli

Knowing stimulus power at any set of points in analysis space provides no information about stimulus power at any other point.



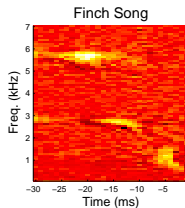
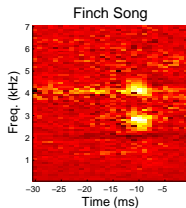
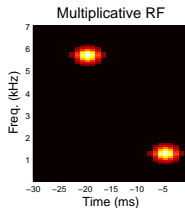
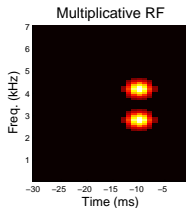
Independence is a property of stimulus *and* analysis space.

Nonlinearity & non-independence distort RF estimates



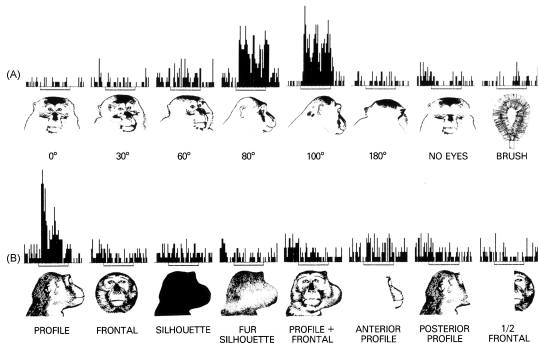
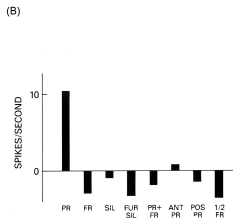
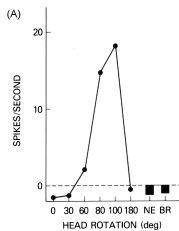
Stimulus may have higher-order correlations in other analysis spaces
— interaction with nonlinearities can produce misleading “receptive fields.”

What about natural sounds?



Usually not independent in any space — so STRFs may not be conservative estimates of receptive fields.

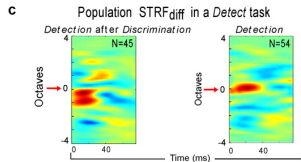
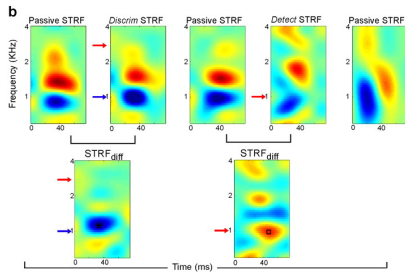
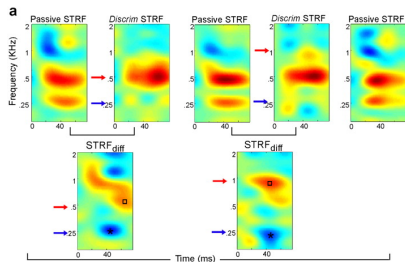
Issues: complex selectivity



a



Issues: adaptation, task-dependence



The “agnostic” coding approach can only take us so far. Eventually, we need solid scientifically (and probably theoretically) motivated hypotheses.