

# PFB inversion

Ningyuan Li

June 2022

## 1 Forward PFB's matrix form

We recap first the basic procedures of an  $r$ -tap PFB. Say, we want to have  $n_c$  channels after the transformation. In a naive FFT procedure, we would need a time stream of length  $2n_c$ , after taking account of the Nyquist criterion<sup>1</sup>. However, in the PFB:

1. We first take a time stream that is of length  $2rn_c$ .
2. We then apply a sinc function of appropriate shape to the original time stream. The sinc is also of length  $2rn_c$ .
3. The sinc-modified time stream is then split into  $r$  pieces of the same length (say,  $M^2$ ). The  $r$  pieces are then stacked together (term-wise added).
4. We finally take the  $M$ -point real-DFT on the stacked signal, which gives us the final answer.

We wish to enumerate some of these steps with linear algebra notation, which makes the inversion procedure more straightforward. Using the prof's notation, we have that:

$$D = FSWd, \tag{1}$$

where  $D$  is the transformed signal,  $d$  is the original time stream, and:

**F:** is the DFT matrix, of size  $M \times M$ . It is **unitary** and hence invertible, which is convenient. The matrix is used in step 4. (Question 1: we are taking the real-DFT instead of the full DFT, which changes the matrix from an invertible one to a non-invertible one (since it would not be square). Is it not a big issue, because the real-DFT is simply the DFT minus some complex conjugates, so we can easily reconstruct the original, and the maths will be the same?)

**S:** is the matrix that takes care of the split-and-stacking. It does what step 3 instructs.

**W:** is the matrix that applies the window. It takes care of step 2.

---

<sup>1</sup>In practice, we take the `rfft` function, which disregards the complex conjugates – and hence having half the number of channels than the original time stream.

<sup>2</sup>Notice that  $M = 2n_c$ .

Writing everything out explicitly, we have that<sup>3</sup>:

$$\underbrace{\begin{bmatrix} D_1 \\ \vdots \\ D_{n_c} \end{bmatrix}}_D = \underbrace{\begin{bmatrix} F_{1,1} & \cdots & F_{1,2n_c} \\ \vdots & \ddots & \vdots \\ F_{n_c,1} & \cdots & F_{n_c,2n_c} \end{bmatrix}}_F \underbrace{\begin{bmatrix} \text{id}_{2n_c} & \cdots & \text{id}_{2n_c} \\ \vdots & \ddots & \vdots \\ \text{id}_{2n_c} & \cdots & \text{id}_{2n_c} \end{bmatrix}}_{\substack{r \\ S}} \underbrace{\begin{bmatrix} W_{1,1} & 0 & \cdots & 0 \\ 0 & W_{2,2} & \ddots & \vdots \\ \vdots & \ddots & \ddots & 0 \\ 0 & \cdots & 0 & W_{2rn_c,2rn_c} \end{bmatrix}}_W \underbrace{\begin{bmatrix} d_1 \\ \vdots \\ d_{2rn_c} \end{bmatrix}}_d$$

It is also of future interest to write the exact form of  $SW$ <sup>4</sup>. Notice that, since  $S$  is a size  $2n_c \times 2rn_c$  matrix, the product would not actually be square, and hence not invertible. Some techniques shall be explored in the future to undo its effect.

$$SW = \begin{bmatrix} W_1 & \cdots & 0 & W_{2n_c+1} & \cdots & 0 & \cdots & W_{2(r-1)n_c+1} & \cdots & 0 \\ \vdots & \ddots & \vdots & \vdots & \ddots & \vdots & \cdots & \vdots & \ddots & \vdots \\ 0 & \cdots & W_{2n_c} & 0 & \cdots & W_{4n_c} & \cdots & 0 & \cdots & W_{2rn_c} \end{bmatrix}, \quad (2)$$

as we can see,  $SW$  is in the form of  $r$  size  $2n_c$  diagonal matrices horizontally aligned.

It may also be interesting to observe that each row of  $SW$  only acts upon samples in  $d$  whose indices are  $2n_c$  apart. This will become important later on.

(Question 2: in the first paragraph of section 3 in your notes, I'm not quite sure what the procedure is that gives a Toeplitz matrix in the end. Is it the result of taking the IDFT (so multiplying  $F^{-1}$  again on  $SW$ )? Or is it something else? Moreover, is the 'decoupled chunk' that you are referring to the horizontal, diagonal blocks in  $SW$ ?)

## 2 PFB of long time stream

Say, we want to compute the PFB of a very long time stream, with the parameters  $r$  and  $n_c$  predetermined. A normal FFT would return  $n_c$  channels from a time stream of length  $2n_c$ , and the PFB simply returns the same number of channels from  $2rn_c$  samples. As such, in order to be able to return the same amount of processed data as ordinary FFT (Question 3: is this the reason why we shift by  $2n_c$ ?), we also shift forward in time by  $2n_c$  samples after every instance of PFB.

To illustrate what this means more concretely, say, we have the original (very long time stream)  $d = [d_1 \ \cdots \ d_{2n_c} \ d_{2n_c+1} \ \cdots \ d_{4n_c} \ \cdots]^T$ , and we wish to compute an  $r$ -tap<sup>5</sup> PFB that outputs  $n_c$  channels per block. If we stack the outputs in a matrix, we would obtain:

$$\mathcal{D}^T = \begin{bmatrix} D_1 & \cdots & D_{n_c} \\ D_{n_c+1} & \cdots & D_{2n_c} \\ \vdots & \ddots & \vdots \end{bmatrix},$$

<sup>3</sup> $\text{id}_s$  stands for identity matrix of size  $s$ .

<sup>4</sup>In below discussion, elements  $W_{i,j} \in W$  will be denoted more compactly as  $W_i$ , since  $i = j \ \forall i, j$ .

<sup>5</sup>Note that the tap number does not impact the sizes of inputs and outputs.

## 2 PFB of long time stream

where the first row consists of the outputs from a PFB of  $d_1$  to  $d_{2rn_c}$ , the second row consists of the outputs from the PFB of  $d_{2n_c}$  to  $d_{2(r+1)n_c}$ , and so on.

More ambitiously, we can try to condense the operation into one single matrix equation,  $\mathcal{D} = \mathcal{F}\mathcal{S}\mathcal{W}d$ .  $\mathcal{F}$  is easily interpreted as the operator that applies the matrix  $F$  (cf. Eq. (1)) to each block of the output, which leaves  $\mathcal{S}\mathcal{W}$  to be the more complicated part. (Question 4: the matrix equation does not actually provide the output  $\mathcal{D}$  as a stacked matrix, since the dimensions of matrix multiplication dictates otherwise. I don't think it's a huge problem since the output of  $\mathcal{S}\mathcal{W}d$  can easily be sliced-and-stacked to produce a stacked matrix; am I wrong to consider it this way?)

To figure out the form that the matrix takes, we notice that the matrix is necessarily some stacked form of Eq. (2), since we are still foundationally doing PFB's to each block. Moreover, each  $SW$  block in the bigger  $\mathcal{S}\mathcal{W}$  matrix must be stacked vertically, in order to make sure that the output in  $\mathcal{D}$  is stacked block-by-block. The other key characteristic of  $\mathcal{S}\mathcal{W}$  is that each  $SW$  block is aligned vertically in a staggered fashion, which is for the purpose of shifting forward by  $2n_c$  samples in  $d$  for instance of PFB computations.

With those in mind, we can deduce the form of  $\mathcal{S}\mathcal{W}$ , with  $SW$  blocks as components:

$$\mathcal{S}\mathcal{W} = \begin{bmatrix} SW & 0 & \cdots & \cdots & 0 \\ \xleftrightarrow[2n_c]{} & SW & 0 & \ddots & 0 \\ \xleftrightarrow[2n_c]{} & \xleftrightarrow[2n_c]{} & SW & \cdots & 0 \\ \cdots & \cdots & \cdots & \cdots & \cdots \end{bmatrix}, \quad (3)$$

note that contrary to this rather poor portrayal, the  $SW$  block are not vertically disjoint from each other; rather, they are staggered, with each block moving  $2n_c$  entries to the right.

*Example 2.1. We now try to illustrate our method with a concrete example. Say, we have the input  $d$  with  $|d| = 12$ , with each individual sample denoted as the corresponding letter of the English alphabet. We would also want to impose a PFB scheme that is 3-tap, with  $2n_c = 3$ <sup>6</sup>. As a result, our matrix equation looks something like this:*

$$\begin{bmatrix} W_1 & & & W_4 & & & W_7 & & & & & \\ & W_2 & & & W_5 & & & W_8 & & & & \\ & & W_3 & & & W_6 & & & W_9 & & & \\ & & & W_1 & & & W_4 & & & W_7 & & \\ & & & & W_2 & & & W_5 & & & W_8 & \\ & & & & & W_3 & & & W_6 & & & W_9 \end{bmatrix} \cdot \begin{bmatrix} a \\ b \\ c \\ \vdots \\ k \\ l \end{bmatrix},$$

*with the empty spaces in the matrix filled with zeros. Now, if we compute the output, we would*

---

<sup>6</sup>Yes, this is not ideal. But for the purpose of illustration it should suffice, as we are not dealing with the final FFT step here.

obtain the following vector:

$$\mathbf{v} = \begin{bmatrix} aW_1 + dW_4 + gW_7 \\ bW_2 + eW_5 + hW_8 \\ cW_3 + fW_6 + iW_9 \\ dW_1 + gW_4 + jW_7 \\ eW_2 + hW_5 + kW_8 \\ fW_3 + iW_6 + lW_9 \end{bmatrix};$$

easily, we could do some basic manipulation to rearrange  $\mathbf{v}$  into the form that we want. Recall that we want each row to be independently ready for the ensuing DFT operation. As a result, we simply have to arrange every three entries in  $\mathbf{v}$  in a row, and we have our desired result. To spell it out, we want:

$$\mathcal{S}\mathcal{W}\mathbf{d} = \begin{bmatrix} v_1 & v_2 & v_3 \\ v_4 & v_5 & v_6 \end{bmatrix},$$

where the notation  $v_k$  means the  $k$ -th entry of  $\mathbf{v}$ . We would get our end result by simply computing the DFT of each row. ■

Now, if we inspect the form of Eq. (3) in combination with Eq. (2), we can see that entries that are  $2n_c$  indices apart in the original time stream  $\mathbf{d}$  can actually be seen as equivalent classes that each forms one decoupled problem. In other words, within  $\mathbf{d}$ , there are  $2n_c$  separate problems, each is its own matrix equation, with the matrix being certain entries from Eq. (3).

With that in mind, we shift our focus to those said matrix operations. The principle observation to make is that each decoupled problem out of the  $2n_c$  total problems is operated by a **Toeplitz matrix**, with  $r$  diagonals for an  $r$ -tap PFB. To demonstrate it more concretely, consider Exp. 2.1 again. In that particular problem, we have three decoupled problems out of the 12 samples in  $\mathbf{d}$ :  $\mathbf{d}_1 = [a \ d \ g \ j]^T$ ,  $\mathbf{d}_2 = [b \ e \ h \ k]^T$ ,  $\mathbf{d}_3 = [c \ f \ i \ l]^T$ . We take  $\mathbf{d}_1$  as an example; by extracting the entries from the preceding  $\mathcal{S}\mathcal{W}$  that operates on this four-long vector, we arrive at the following matrix equation:

$$\text{output} = T_1 \mathbf{d}_1 = \begin{bmatrix} W_1 & W_4 & W_7 & 0 \\ 0 & W_1 & W_4 & W_7 \end{bmatrix} \begin{bmatrix} a \\ d \\ g \\ j \end{bmatrix}. \quad (4)$$

The matrix  $T_1$  clearly takes care of all the operations that is responsible for the subproblem.

As we have foreshadowed, every single instance of  $\mathbf{d}_k : k \in \{1, \dots, 2n_c\}$  is operated by a corresponding Toeplitz matrix  $T_k$ . This situation is interesting, in that the multiplication of a Toeplitz matrix with a vector actually is another way of computing convolutions. We can write out an example just to convince ourselves:

*Example 2.2. We wish to show that a convolution between two vectors can be represented easily using a Toeplitz matrix. For simplicity, we stick to vectors of small cardinality for the example.*

Consider vector  $\mathbf{a}$  of length three and vector  $\mathbf{b}$  of length four. We wish to compute their convolution  $\mathbf{a} * \mathbf{b} = \mathbf{c}$ , clearly of length  $3 + 4 - 1 = 6$ . Using the definition of the convolution operation, we can compute each entry of  $\mathbf{c}$  rather easily:

$$\mathbf{c} = \begin{bmatrix} a_1 b_1 \\ a_1 b_2 + b_1 a_2 \\ a_1 b_3 + b_2 a_2 + a_3 b_1 \\ a_1 b_4 + a_2 b_3 + a_3 b_2 \\ a_2 b_4 + a_3 b_2 \\ a_3 b_4 \end{bmatrix}.$$

We can also do the same with a matrix, however:

$$\mathbf{c} = T_b \mathbf{a} = \begin{bmatrix} b_1 & & & \\ b_2 & b_1 & & \\ b_3 & b_2 & b_1 & \\ b_4 & b_3 & b_2 & \\ & b_4 & b_3 & \\ & & b_4 & \end{bmatrix} \cdot \begin{bmatrix} a_1 \\ a_2 \\ a_3 \end{bmatrix},$$

where empty entries are zeros. As we can see,  $T_b$  is indeed a Toeplitz matrix. ■

The Toeplitz-ness will come in handy when we try to invert the problem.

### 3 Inverting the PFB

The big problem with convolutions is that they are not always invertible (like what's happening in Exp. 2.2). As evidenced by Eq. (4), after manipulating the input data with the PFB, the output vector is smaller in dimension compared to the input. This is a problem, since we could clearly not construct a bijection – and hence no inversion – between pre- and post-PFB data. Nevertheless, we could still try and do some manipulations that allows us to ‘invert’ the PFB, despite always imperfectly.

(From now on everything I do is a wild guess at attempting to figure out what you meant in the memo. I’m guessing on paper for the moment to work out some maths since I’m stuck currently.)

We first turn our gaze into the size of each Toeplitz matrix, and its relationship with the parameters of the PFB. Again, we want an  $r$ -tap PFB with  $n_c$  channels of output. We also have an input signal of size  $s$ <sup>7</sup>. Naturally, since we divide the data set to  $2n_c$  decoupled problems, each Toeplitz matrix must have  $\frac{s}{2n_c}$  columns – which is also the size of each sub-problem.

Knowing the above information, we could also compute the number of rows of the matrix. Since the PFB is  $r$ -tap, the matrix must have  $\frac{s}{2n_c} - r + 1$  rows. This means that

---

<sup>7</sup>We always assume divisibility here, since the guideline to a long data set is to delete the remainder.

as the size of each decoupled problem increases, the matrix becomes more and more square. This is good news, since in a certain sense it enables us to ‘approximate’ the Toeplitz matrix as being square. And, to make our life easier, we complete the missing  $r - 1$  rows such that the resultant matrix is **circulant** – meaning that each row is the right-cyclical permutation of the preceding row by one entry. Doing this lends many interesting properties of the circulant matrix to our disposal, and it greatly aids our inversion procedure.

We still have a few pieces of the puzzle left before the actual inverse. Firstly, consider the original equation for the decoupled problems:

$$\text{output} = T_k d_k,$$

where  $d_k$  is the  $k$ -th problem out of the  $2n_c$  ones, and  $T_k$  is its corresponding Toeplitz matrix. We have to recognise that the difference between the number of rows and columns of  $T_k$  means that after we complete it to be circulant, i.e.:

$$\text{output}' = C_k d_k,$$

where  $C_k$  is the completed circulant matrix, the last few rows of  $C_k$  would actually yield something nonsensical for the last few entries of  $\text{output}'$ . Now, in the context of the PFB, we have hundreds (or millions) of data points to process in each instance of the matrix equation. As a result, in the grand scheme of things, it is not a huge deal to have a few entries there. Because of this, we simply put the first few entries of the output, and fill those entries back in the end such that the dimension of the output matches that of the input.

What is now left for us to do is simply the procedure of inverting each  $C_k$ . We have mentioned before that multiplying by a Toeplitz matrix is equivalent to a convolution operation, and since circulant matrices are just special types of Toeplitz matrices, this still holds true here. Interestingly, this fact actually hints to the invertibility of circulant matrices. By the convolution theorem and some other maths (such as solving difference equations and guessing for a few *Ansätze*), we can derive that the eigenvectors of the (or any, in fact) circulant matrix consists of entries in order from the zeroth power of the  $k$ -th power of the  $n$ -th root of unity to the  $(n - 1)$ -th power, assuming that the circulant matrix in question is an  $n \times n$  matrix, and the entries of the first row are indexed with  $k$ .

We can try to express the same idea but more intelligibly. Say, we have that  $C$  is circulant of size  $n$ , and we denote the first row of  $C$  as  $c = [c_0 \cdots c_k \cdots c_{n-1}] : k \in [0, n - 1]$ . Then, the  $k$ -th eigenvector of  $C$  turns out to be:

$$v = \begin{bmatrix} \omega_n^{0k} \\ \omega_n^{1k} \\ \vdots \\ \omega_n^{(n-1)k} \end{bmatrix},$$

where  $\omega_n = e^{\frac{2\pi i}{n}}$ . The corresponding eigenvalues can also be derived, as the entries of the DFT of  $c$ . This fact is remarkable, in that it shows that any circulant matrix can be diagonalised –  $C = U\Lambda U^{-1}$ . The other interesting property of this procedure is that after some re-normalisation, the eigenvectors stacked together in order as column vectors actually gives us the DFT matrix,  $F$  (far cry from the first section). More precisely, we have that  $U = \frac{1}{\sqrt{n}}F_n$ . As a result, in order to invert the equation:

$$\text{output}' = C_k \mathbf{d}_k,$$

we simply have to do some algebraic manipulations:

$$\begin{aligned} \text{output}' &= \frac{1}{n} F \Lambda_k F^{-1} \mathbf{d}_k, \\ n F \Lambda_k^{-1} F^{-1} \text{output}' &= \mathbf{d}_k. \end{aligned}$$

Disregarding the constant  $n$ , what is left is simply a simple three-step procedure:

1. take the IDFT of the output';
2. take the DFT of the first row of  $C_k$ , then divide each entry of step 1 with the result of the DFT of matching indices;
3. take the DFT of the result from step 2.

(Question 5: what I have derived seems similar to what you have described on your memo, but with one key difference: in your memo, it seems that your inversion procedure is to take the DFT, divide by the DFT'ed first row, then take the IDFT. However, from the maths that I had tried to work out, the first and third step seems to be inverted. Did I do something wrong along the way in my maths?)