

# Fast generalized DFTs for all finite groups

Chris Umans\*  
Caltech

January 10, 2019

## Abstract

For any finite group  $G$ , we give an arithmetic algorithm to compute generalized Discrete Fourier Transforms (DFTs) with respect to  $G$ , using  $O(|G|^{\omega/2+\epsilon})$  operations, for any  $\epsilon > 0$ . Here,  $\omega$  is the exponent of matrix multiplication.

## 1 Introduction

For a finite group  $G$ , let  $\text{Irr}(G)$  denote a complete set of irreducible representations of  $G$ . A *generalized DFT with respect to  $G$*  is a map from a group algebra element  $\alpha \in \mathbb{C}[G]$  (which is a vector of  $|G|$  complex numbers), to the following linear combination of irreducible representations:

$$\sum_{g \in G} \alpha_g \bigoplus_{\rho \in \text{Irr}(G)} \rho(g).$$

It is unique once one fixes a basis for each  $\rho$ ; one usually seeks algorithms that work for arbitrary chosen bases. We typically speak of the complexity of computing this map in the (non-uniform) arithmetic circuit model and do not concern ourselves with *finding* the irreducible representations. The trivial algorithm thus requires  $O(|G|^2)$  operations, since we are summing  $|G|$  block-diagonal matrices, each with  $|G|$  entries in the blocks.

Fast algorithms for the DFT with respect to cyclic groups are well-known and are attributed to Cooley and Tukey in 1965 [CT65], although the ideas likely date to Gauss. Beth in 1984 [Bet84], together with Clausen [Cla89], initiated the study of generalized DFTs, the “generalized” terminology signalling that the underlying group may be any group. A central goal since that time has been to obtain fast algorithms for generalized DFTs with respect to arbitrary underlying groups. One may hope for “nearly-linear” time algorithms, meaning that they use a number of operations that is upper-bounded by  $c_\epsilon |G|^{1+\epsilon}$  for universal constants  $c_\epsilon$  and arbitrary  $\epsilon > 0$ . Such “exponent one” algorithms are known for certain families of groups: abelian groups, supersolvable groups [Bau91], and symmetric and alternating groups [Cla89]. Algorithms for generalized DFTs often find themselves manipulating matrices, so it is not surprising that they require a number of operations that depends on  $\omega$ , the exponent of matrix multiplication. Thus we view algorithms that achieve exponent one conditioned on  $\omega = 2$  as being “nearly as good” as unconditional exponent one algorithms. Such algorithms are known for solvable groups [Bet84, CB93], and with the recent breakthrough of [HU18a], for linear groups; these algorithms achieve exponent  $\omega/2$ .

In this paper we realize the main goal of the area, obtaining exponent  $\omega/2$  for all finite groups  $G$ . The previous best exponent that applies to all finite groups was obtained by [HU18a]; it depends in a somewhat complicated way on  $\omega$ , but it is at best  $\sqrt{2}$  (when  $\omega = 2$ ). Before that, the best known exponent was  $1 + \omega/4$  (which is at best  $3/2$  when  $\omega = 2$ ), and this dates back to the original work of Beth and Clausen.

---

\*Supported by NSF grant CCF-1815607 and a Simons Foundation Investigator grant.

## 1.1 Past and related work

A good description of past work in this area can be found in Section 13.5 of [BCS97]. The first algorithm generalizing beyond the abelian case is due to Beth in 1984 [Bet84]; this algorithm is described in Section 3.1 in a form often credited jointly to Beth and Clausen. Three other milestones are the  $O(|G|\log|G|)$  algorithm for supersolvable groups due to Baum [Bau91], the  $O(|G|\log^3|G|)$  algorithm for the symmetric group due to Clausen [Cla89] (see also [Mas98] for a recent improvement), and the  $O(|G|^{\omega/2+\epsilon})$  algorithms for linear groups obtained by Hsu and Umans, which is described in Section 3.2. Wreath products were studied by Rockmore [Roc95] who obtained exponent one algorithms in certain cases.

In the 1990s, Maslen, Rockmore, and coauthors developed the so-called “separation of variables” approach [MR97a], which relies on non-trivial decompositions along chains of subgroups via *Bratteli diagrams* and detailed knowledge of the representation theory of the underlying groups. There is a rather large body of literature on this approach and it has been applied to a wide variety of group algebras and more general algebraic objects. For a fuller description of this approach and the results obtained, the reader is referred to the surveys [MR97b, Roc02], and the most recent paper in this line of work [MRW16].

## 2 Preliminaries

Throughout this paper we will use the phrase

“generalized DFTs with respect to  $G$  can be computed using  $O(|G|^{\alpha+\epsilon})$  operations, for all  $\epsilon > 0$ ”

where  $G$  is a finite group and  $\alpha \geq 1$  is a real number. We mean by this that there are *universal* constants  $c_\epsilon$  independent of the group  $G$  under consideration so that for each  $\epsilon > 0$ , the operation count is at most  $c_\epsilon|G|^{\alpha+\epsilon}$ . Such an algorithm will be referred to as an “exponent  $\alpha$ ” algorithm. This comports with the precise definition of the exponent of matrix multiplication,  $\omega$ : that there are universal constants  $b_\epsilon$  for which  $n \times n$  matrix multiplication can be performed using at most  $b_\epsilon n^{\omega+\epsilon}$  operations, for each  $\epsilon > 0$ . Indeed we will often report our algorithms’ operation counts in terms of  $\omega$ . In such cases matrix multiplication is always used as a black box, so, for example, an operation count of  $O(|G|^{\omega/2})$  should be interpreted to mean: if one uses a fast matrix multiplication algorithm with exponent  $\alpha$  (which may range from 2 to 3), then the operation count is  $O(|G|^{\alpha/2})$ . In particular, in real implementations, one might well use standard matrix multiplication and plug in 3 for  $\omega$  in the operation count bound.

We use  $\text{Irr}(G)$  to denote the complete set of irreducible representations of  $G$  being used for the DFT at hand. In the presentation to follow, we assume the underlying field is  $\mathbb{C}$ ; however our algorithms work over any field  $\mathbb{F}_{p^k}$  whose characteristic  $p$  does not divide the order of the group, and for which  $k$  is sufficiently large for  $\mathbb{F}_{p^k}$  to represent a complete set of irreducibles.

We use  $I_n$  to denote the  $n \times n$  identity matrix. The following is an important general observation (see, e.g., Lemma 4.3.1 in [HJ91]):

**Proposition 2.1.** *If  $A$  is an  $n_1 \times n_2$  matrix,  $B$  is an  $n_2 \times n_3$  matrix, and  $C$  is an  $n_3 \times n_4$  matrix, then the entries of the product matrix  $ABC$  are exactly the entries of the vector obtained by multiplying  $A \otimes C^T$  (which is an  $n_1 n_4 \times n_2 n_3$  matrix) by  $B$  viewed as an  $n_2 n_3$ -vector, and denoted  $\text{vec}(B)$ .*

### 2.1 Basic representation theory

A *representation* of group  $G$  is a homomorphism  $\rho$  from  $G$  into the group of invertible  $d \times d$  matrices. Representation  $\rho$  naturally specifies an action of  $G$  on  $\mathbb{C}^d$ ; representation  $\rho$  is thus said to have *dimension*  $\dim(\rho) = d$ . A representation is *irreducible* if the action on  $\mathbb{C}^d$  has no  $G$ -invariant subspace. Two representations of the same dimension  $d$ ,  $\rho_1$  and  $\rho_2$ , are *equivalent* (written  $\rho_1 \cong \rho_2$ ) if they are the same up to a change of basis; i.e.,  $\rho_1(g) = T\rho_2(g)T^{-1}$  for some invertible  $d \times d$  matrix  $T$ . The classical Maschke’s Theorem implies that every representation

$\rho_0$  of  $G$  breaks up into the direct sum of irreducible representations; i.e. there is an invertible matrix  $T$  and a multiset  $S \subseteq \text{Irr}(G)$ , for which

$$T\rho_0(g)T^{-1} = \bigoplus_{\rho \in S} \rho(g).$$

Given a subgroup  $H \subseteq G$  one can obtain from any representation  $\rho \in \text{Irr}(G)$  a representation  $\text{Res}_H^G(\rho)$  (the *restriction* of  $\rho$  to  $H$ ), which is a representation of  $H$ , simply by restricting the domain of  $\rho$  to  $H$ . One can also obtain from any representation  $\sigma \in \text{Irr}(H)$ , a representation of  $G$  called the *induced* representation  $\text{Ind}_H^G(\sigma)$ , which has dimension  $\dim(\sigma)|G|/|H|$ . We will not need to work directly with induced representations, but we will use a fundamental fact called *Frobenius reciprocity*. Given  $\rho \in \text{Irr}(G)$  and  $\sigma \in \text{Irr}(H)$ , Frobenius reciprocity states that the number of times  $\sigma$  appears in the restriction  $\text{Res}_H^G(\rho)$  equals the number of times  $\rho$  appears in the induced representation  $\text{Ind}_H^G(\sigma)$ .

A basic fact is that  $\sum_{\rho \in \text{Irr}(G)} \dim(\rho)^2 = |G|$ , which implies that for all  $\rho \in \text{Irr}(G)$ , we have  $\dim(\rho) \leq |G|^{1/2}$ . This can be used to prove the following inequality, which we use repeatedly:

**Proposition 2.2.** *For any real number  $\alpha \geq 2$ , we have*

$$\sum_{\rho \in \text{Irr}(G)} \dim(\rho)^\alpha \leq |G|^{\alpha/2}.$$

*Proof.* Set  $\rho_{\max}$  to be an irrep of largest dimension. We have

$$\sum_{\rho \in \text{Irr}(G)} \dim(\rho)^\alpha \leq \dim(\rho_{\max})^{\alpha-2} \sum_{\rho \in \text{Irr}(G)} \dim(\rho)^2 = \dim(\rho_{\max})^{\alpha-2} |G| \leq |G|^{\alpha/2},$$

where the last inequality used the fact that  $\dim(\rho_{\max}) \leq |G|^{1/2}$ . □

## 2.2 Basic Clifford theory

Clifford theory describes the way the irreducible representations of a group  $H$  break up when restricted to a *normal* subgroup  $N$ , which is a particularly well-structured and well-understood scenario.

Elements of  $H$  act on the set  $\text{Irr}(N)$  as follows:

$$(h \cdot \lambda)(n) = \lambda(hnh^{-1}),$$

for  $\lambda \in \text{Irr}(N)$ . Let  $\mathcal{O}_1, \dots, \mathcal{O}_\ell$  be the orbits of this  $H$ -action on  $\text{Irr}(N)$ . Clifford theory states for each  $\sigma \in \text{Irr}(H)$ , there is a positive integer  $e_\sigma$  and an index  $i_\sigma$  for which the restriction  $\text{Res}_N^H(\sigma)$  is equivalent to

$$e_\sigma \bigoplus_{\lambda \in \mathcal{O}_{i_\sigma}} \lambda.$$

In particular, this implies that all  $\lambda \in \text{Irr}(N)$  that occur in the restriction have the same dimension,  $d_\sigma$ , and multiplicity,  $e_\sigma$ , and that  $\dim(\sigma) = d_\sigma e_\sigma |\mathcal{O}_{i_\sigma}|$ .

We can also define the following subsets, which partition  $\text{Irr}(H)$ :

$$S_\ell = \{\sigma \in \text{Irr}(H) : \text{the irreps in } \mathcal{O}_\ell \text{ occur in } \sigma\} = \{\sigma \in \text{Irr}(H) : i_\sigma = \ell\}.$$

We will need the following proposition:

**Proposition 2.3.** *For a finite group  $H$  and normal subgroup  $N$ , and sets  $S_\ell$  as defined above, the following holds for each  $\ell$ :*

$$\sum_{\sigma \in S_\ell} \dim(\sigma) e_\sigma / d_\sigma = |H/N|.$$

*Proof.* Fix  $\lambda \in \mathcal{O}_\ell$ , and note that the induced representation  $\text{Ind}_N^H(\lambda)$  has dimension  $\dim(\lambda)|H/N|$ . Let  $m_{\sigma,\lambda}$  be the number of times  $\sigma \in \text{Irr}(H)$  occurs in  $\text{Ind}_N^H(\lambda)$ . Then we have

$$\sum_{\sigma \in \text{Irr}(H)} \dim(\sigma)m_{\sigma,\lambda} = \dim(\lambda)|H/N|.$$

By Frobenius reciprocity,  $m_{\sigma,\lambda}$  equals the number times  $\lambda$  occurs in  $\text{Res}_N^H(\sigma)$ . Thus the summand  $\dim(\sigma)m_{\sigma,\lambda}$  equals  $\dim(\sigma)e_{\sigma,\lambda}$ , whenever  $m_{\sigma,\lambda} \neq 0$  (and zero otherwise). The proposition follows.  $\square$

### 2.3 Generalized DFTs and inverse generalized DFTs

We assume by default that we are computing generalized DFTs with respect to an arbitrary chosen basis for each  $\rho \in \text{Irr}(G)$ . Sometimes we need to refer to the special basis in the following definition:

**Definition 2.4.** *Let  $H$  be a subgroup of  $G$ . An  $H$ -adapted basis is a basis for each  $\rho \in \text{Irr}(G)$ , so that the restriction of  $\rho$  to  $H$  respects the direct sum decomposition into irreps of  $H$ .*

In concrete terms, this implies that for each  $\rho \in \text{Irr}(G)$ , while for general  $g \in G$ ,  $\rho(g)$  is a  $\dim(\rho) \times \dim(\rho)$  matrix, for  $g \in H$ ,  $\rho(g)$  is a block-diagonal matrix with block sizes coming from the set  $\{\dim(\sigma) : \sigma \in \text{Irr}(H)\}$ . An  $H$ -adapted basis always exists.

A general trick that we will rely on is that if one can compute generalized DFTs with respect to  $G$  for an input  $\alpha$  supported on a subset  $S \subseteq G$ , then with an additional multiplicative factor of roughly  $|G|/|S|$ , one can compute generalized DFTs with respect to  $G$ .

**Theorem 2.5.** *Fix a finite group  $G$  and a subset  $S \subseteq G$ , and suppose a generalized DFT with respect to  $G$  can be computed in  $m$  operations, for inputs  $\alpha$  supported on  $S$ . Then generalized DFTs with respect to  $G$  can be computed using*

$$O(m + |G|^{\omega/2+\epsilon}) \cdot \frac{|G| \log |G|}{|S|}$$

operations, for any  $\epsilon > 0$ .

*Proof.* First observe that by multiplying by  $\bigoplus_{\rho \in \text{Irr}(G)} \rho(g)$  we can compute a generalized DFT supported on  $Sg$ , for an additive extra cost of

$$\sum_{\rho \in \text{Irr}(G)} O(\dim(\rho)^{\omega+\epsilon})$$

operations, for all  $\epsilon > 0$ , and by applying Proposition 2.2 with  $\alpha = \omega + \epsilon$  this is at most  $O(|G|^{\omega/2+\epsilon})$ . A probabilistic argument shows that  $|G| \log |G|/|S|$  different translates  $g$  suffice to cover  $G$ , so we need only repeat the DFT supported on  $Sg$  translated by each such  $g$ , and sum the resulting DFTs.  $\square$

The *inverse generalized DFT* maps a collection of matrices  $M^\sigma \in \mathbb{C}^{\dim(\sigma) \times \dim(\sigma)}$ , one for each  $\sigma \in \text{Irr}(G)$ , to the vector  $\alpha$  for which

$$\sum_{g \in G} \alpha_g \bigoplus_{\sigma \in \text{Irr}(G)} \rho(G) = \bigoplus_{\sigma \in \text{Irr}(G)} M^\sigma.$$

In the arithmetic circuit model, the inverse DFT can be computed efficiently if the DFT can:

**Theorem 2.6** (Baum, Clausen; Cor. 13.40 in [BCS97]). *Fix a generalized DFT with respect to finite group  $G$  and suppose it can be computed in  $m$  operations. Then the inverse DFT with respect to  $G$  (and the same basis), can be computed in at most  $m + |G|$  operations.*

### 3 General strategy: reduction to subgroups

One way to organize the main algorithmic ideas in the quest for a fast DFT for all finite groups is according to the subgroup structure they exploit. The algorithms themselves are recursive, with the main content of the algorithm being the reduction to smaller instances: DFTs over subgroups of the original group. When aiming for generalized DFTs for all finite groups, such a reduction is paired with a group-theoretic structural result, which guarantees the existence of certain subgroups that are used by the reduction.

In the exposition below, it is helpful to assume that  $\omega = 2$  and seek an “exponent 1” algorithm under this assumption (in general, the exponent achieved will be a function of  $\omega$ , and in our main result this function is  $\omega/2$ ). By the term *overhead* we mean the extra multiplicative factor in the operation count of the reduction, beyond the nearly-linear operation count that would be necessary for an exponent 1 algorithm.

#### 3.1 The single subgroup reduction

The seminal Beth-Clausen algorithm reduces computing a DFT over a group  $G$  to computing several DFTs over a subgroup  $H$  of  $G$ . We call this the “single subgroup reduction”. Roughly speaking, the overhead in this reduction is proportional to the index of  $H$  in  $G$ . The companion structural result is Lev’s Theorem [Lev92], which shows that every finite group  $G$  (except cyclic of prime order which can be handled separately) has a subgroup of order at least  $\sqrt{|G|}$  (and this is tight, hence the overhead is  $\sqrt{|G|}$  in the worst case). As noted in the introduction, this reduction together with Lev’s Theorem implies exponent  $3/2$  (assuming  $\omega = 2$ ) for all finite groups.

Here is a more detailed description, together with results we will need later. Let  $H$  be a subgroup of  $G$  and let  $X$  be a set of distinct coset representatives. We first compute several  $H$ -DFTs, one for each  $x \in X$ :

$$s_x = \sum_{h \in H} \alpha_{hx} \bigoplus_{\sigma \in \text{Irr}(H)} \sigma(h)$$

and by using an  $H$ -adapted basis (Definition 2.4), we can lift each  $s_x$  to

$$\bar{s}_x = \sum_{h \in H} \alpha_{hx} \bigoplus_{\rho \in \text{Irr}(G)} \rho(h)$$

by just copying entries (which is free of cost in the arithmetic model). Then to complete the DFT we need to compute

$$\sum_{x \in X} \bar{s}_x \bigoplus_{\rho \in \text{Irr}(G)} \rho(x).$$

Generically, this final computation requires an overhead proportional to  $|X| = [G : H]$ , even when just considering the outermost summation. See Corollary 4 in [HU18b] for the details to complete this sketch, yielding the following:

**Theorem 3.1** (single subgroup reduction). *Let  $G$  be a finite group and let  $H$  be a subgroup. Then we can compute a generalized DFT with respect to  $G$  at a cost of  $[G : H]$  many  $H$ -DFTs plus  $O([G : H]|G|^{\omega/2+\epsilon})$  operations, for all  $\epsilon > 0$ .*

In the special case that  $H$  is normal in  $G$  and  $G/H$  is cyclic of prime order, the overhead of  $[G : H]$  can be avoided, by using knowledge about the way representations  $\sigma \in \text{Irr}(H)$  extend to  $\rho \in \text{Irr}(G)$ . This insight is the basis for the Beth-Clausen algorithm for solvable groups. We need it here to handle the case of  $G/H$  cyclic of prime order, which is the single exceptional case not handled by our main reduction. The following theorem can be inferred from the proof of Theorem 7.7 in Clausen and Baum’s monograph [CB93]:

**Theorem 3.2** (Clausen, Baum [CB93]). *Let  $H$  be a normal subgroup of  $G$  with prime index  $p$ . We can compute a generalized DFT with respect to  $G$  and an  $H$ -adapted basis, at a cost of  $p$  many  $H$ -DFTs plus*

$$O(p \log p) \cdot \sum_{\sigma \in \text{Irr}(H)} \dim(\sigma)^{\omega+\epsilon}$$

operations, for all  $\epsilon > 0$ .

For our purposes the following slightly coarser bound suffices, which accomodates an arbitrary basis change (hence obviating the need for an  $H$ -adapted basis):

**Corollary 3.3.** *Let  $H$  be a normal subgroup of  $G$  with prime index  $p$ . Generalized DFTs with respect to  $G$  can be computed at a cost of  $p$  many  $H$ -DFTs plus  $O(|G|^{\omega/2+\epsilon})$  operations, for all  $\epsilon > 0$ .*

*Proof.* Applying Proposition 2.2 to Theorem 3.2 with  $\alpha = \omega + \epsilon$  yields an operation count of  $O(p \log p)|H|^{\omega+\epsilon/2}$ , which is at most  $O(|G|^{\omega/2+\epsilon})$ . Performing an arbitrary basis change costs

$$\sum_{\rho \in \text{Irr}(G)} O(\dim(\rho)^{\omega+\epsilon})$$

operations which is again at most  $O(|G|^{\omega/2+\epsilon})$  by Proposition 2.2.  $\square$

### 3.2 The double subgroup reduction

Recently, Hsu and Umans proposed a “double subgroup reduction” [HU18a] which reduces computing a DFT over a group  $G$  to computing several DFTs over two subgroups,  $H$  and  $K$ . This reduction is especially effective for linear groups (see [HU18a]). Roughly speaking, the overhead in this reduction is proportional to  $|G|/|HK|$  and  $|H \cap K|$ . The companion structural result shows that every finite group  $G$  (except  $p$ -groups which can be handled separately) has two proper subgroups  $H$  and  $K$  for which  $|G|/|HK|$  is negligible. However,  $|H \cap K|$  might still be large, which is the one thing standing in the way of deriving an “exponent  $\omega/2$ ” algorithm from this reduction.

To illustrate the bottleneck in this reduction, we describe it in more detail. Let  $H, K$  be subgroups of  $G$  and assume  $|G|/|HK|$  is negligible. We first compute an intermediate representation

$$\sum_{g=hk \in HK} \alpha_g \bigoplus_{\substack{\sigma \in \text{Irr}(H) \\ \tau \in \text{Irr}(K)}} \sigma(h) \otimes \tau(k)$$

in two steps (and then lift it to a  $G$ -DFT). The first of the two steps is to compute at most  $[G : H]$  many  $H$ -DFTs, yielding, for each  $k \in K' \subseteq K$  (where  $K'$  is a set of distinct coset representatives of  $H$  in  $G$ ):

$$s_k = \sum_{h \in H} \alpha_{hk} \bigoplus_{\sigma \in \text{Irr}(H)} \sigma(h).$$

The second step is as follows: for each *entry* of the block-diagonal matrix  $s_k$ , we use this entry (as  $k$  varies) as the data for a  $K$ -DFT. There are  $\sum_{\sigma \in \text{Irr}(H)} \dim(\sigma)^2 = |H|$  such entries in general. Thus the second step entails  $|H|$  many  $K$ -DFTs, and this represents the key bottleneck. Note that when  $|G|/|HK|$  is negligible,  $|H||K|$  is approximately  $|G||H \cap K|$ , and this explains the overhead of roughly  $|H \cap K|$  which prevents obtaining an “exponent  $\omega/2$ ” algorithm from this reduction. For completeness we record the main theorem of [HU18b] here:

**Theorem 3.4** (Theorem 12 in [HU18b]). *Let  $G$  be a finite group and let  $H, K$  be subgroups. Then we can compute generalized DFTs with respect to  $G$  at the cost of  $|H|$  many  $K$ -DFTS,  $|K|$  many  $H$ -DFTs, plus*

$$O(|G|^{\omega/2+\epsilon} + (|H||K|)^{\omega/2+\epsilon})$$

operations, all repeated  $O(\frac{|G| \log |G|}{|HK|})$  times, for all  $\epsilon > 0$ .

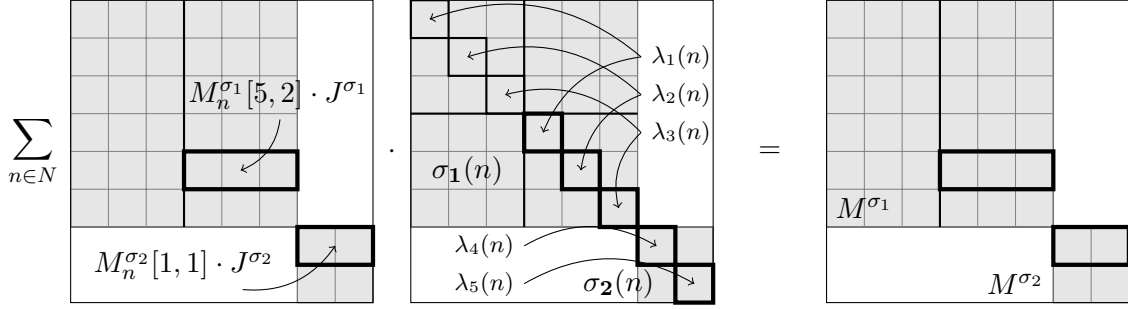


Figure 1: Illustration of the proof of Theorem 3.5. In this example  $\text{Irr}(H) = \{\sigma_1, \sigma_2\}$ ,  $\text{Irr}(N) = \{\lambda_1, \lambda_2, \lambda_3, \lambda_4, \lambda_5\}$ ; the orbits are  $\mathcal{O}_1 = \{\lambda_1, \lambda_2, \lambda_3\}$  and  $\mathcal{O}_2 = \{\lambda_4, \lambda_5\}$ ;  $S_1 = \{\sigma_1\}$  and  $S_2 = \{\sigma_2\}$ ; and the multiplicities are  $e_{\sigma_1} = 2$  and  $e_{\sigma_2} = 1$ . In the figure, we highlight the parts of the matrices that give rise to the system of equations solved with a single inverse  $N$ -DFT, corresponding to the value  $a = f_1(\sigma_1, 5, 2) = f_2(\sigma_2, 1, 1)$ . This inverse  $N$ -DFT with the highlighted blocks of  $M^{\sigma_1}$  and  $M^{\sigma_2}$  as input data yields the scalars  $M_n^{\sigma_1}[5, 2] = M_n^{\sigma_2}[1, 1]$  that satisfy the simultaneous equations.

Our main innovation, described in the next section, is a way to overcome the bottleneck: when  $H \cap K = N$  is a normal subgroup of  $G$ , we are able to rewrite each  $s_k$  as a sum of  $|N|$  matrices with special structure: effectively, there are only  $|H/N|$  many non-zero “entries” for which we need to compute a  $K$ -DFT, and as we will show, this exactly removes the overhead factor.

### 3.3 The triple subgroup reduction

In this section we give our main new result. We devise a “triple subgroup reduction” which reduces computing a DFT over  $G$  to computing several DFTs over two subgroups,  $H$  and  $K$ , and several inverse DFTs over the intersection  $N = H \cap K$ , when  $N$  is normal in  $G$ . Roughly speaking, the overhead is proportional to  $|G|/|HK|$ . The companion structural result (Theorem 3.10) shows that for every finite group  $G$ , if  $N$  is a maximal normal subgroup in  $G$  then (except for the case of  $|G/N|$  cyclic of prime order, which can be handled separately) there exist two proper subgroups  $H$  and  $K$  with  $H \cap K = N$ , such that  $|G|/|HK|$  is negligible. This is the key to the claimed exponent  $\omega/2$  algorithm.

Let  $H$  be a group with normal subgroup  $N$ . The main technical theorem shows how to rewrite the output of an  $H$ -DFT as the sum of  $|N|$  matrices each of which only has “ $|H/N|$  degrees of freedom”. In the following theorem we adopt the notation introduced in Section 2.2.

**Theorem 3.5.** *Let  $H$  be a group and  $N$  a normal subgroup. For every*

$$M = \bigoplus_{\sigma \in \text{Irr}(H)} M^\sigma \in \bigoplus_{\sigma \in \text{Irr}(H)} \mathbb{C}^{\dim(\sigma) \times \dim(\sigma)},$$

*the following holds with respect to an  $N$ -adapted basis: there exist matrices  $M_n^\sigma \in \mathbb{C}^{\dim(\sigma)/d_\sigma \times e_\sigma}$  for which*

$$\sum_{n \in N} (M_n^\sigma \otimes J^\sigma) \cdot \sigma(n) = M^\sigma,$$

*where  $J^\sigma$  is the  $d_\sigma \times \dim(\sigma)/e_\sigma$  matrix  $(I_{d_\sigma} | I_{d_\sigma} | \cdots | I_{d_\sigma})$ . Moreover, given injective functions  $f_\ell$  from  $\{(\sigma, i, j) : \sigma \in S_\ell, i \in [\dim(\sigma)/d_\sigma], j \in [e_\sigma]\}$  to  $[r]$ , the  $M_n^\sigma$  can be taken to satisfy*

$$f_\ell(\sigma, i, j) = f_{\ell'}(\sigma', i', j') \quad \Rightarrow \quad \forall n \quad M_n^\sigma[i, j] = M_n^{\sigma'}[i', j'],$$

*and these matrices  $M_n^\sigma$  can be obtained from  $M$  by computing  $r$  inverse  $N$ -DFTs.*

One should think of the functions  $f_\ell$  as labeling the entries of the  $M_n^\sigma$  matrices for the  $\sigma$  in a given  $S_\ell$ . This labeling is then used to ensure that entries of  $M_n^\sigma$  with  $\sigma \in S_\ell$  and the entries of  $M_n^{\sigma'}$  with  $\sigma' \in S_{\ell'}$  are equal, if they have the same labels. In Section 3.3.1 we will show how to choose this labeling so that the final “lifting” step of our algorithm can be efficiently computed. For now, we note that Proposition 2.3 implies that there *exist* labellings  $f_\ell$  with  $r = |H/N|$ , and indeed our actual choice of  $f_\ell$  in Section 3.3.1 will have  $r = O(|H/N| \log |H/N|)$ , which is not much larger.

*Proof.* Fix  $\sigma \in \text{Irr}(H)$ , and recall that there is a unique  $S_\ell$  containing  $\sigma$ . Since we are using an  $N$ -adapted basis,  $\sigma(n)$  has the form

$$I_{e_\sigma} \otimes \bigoplus_{\lambda \in \mathcal{O}_\ell} \lambda(n),$$

and thus

$$\sum_{n \in N} (M_n^\sigma \otimes J^\sigma) \cdot \sigma(n) = \sum_{n \in N} M_n^\sigma \otimes (\lambda_1(n) |\lambda_2(n)| \cdots |\lambda_{|\mathcal{O}_\ell|}(n)) \quad (1)$$

where  $\lambda_1, \dots, \lambda_{|\mathcal{O}_\ell|}$  is an enumeration of  $\mathcal{O}_\ell$ . Since these are pairwise inequivalent irreps, the span of

$$\{(\lambda_1(n) |\lambda_2(n)| \cdots |\lambda_{|\mathcal{O}_\ell|}(n)) : n \in N\}$$

is the full matrix algebra  $\mathbb{C}^{d_\sigma \times \dim(\sigma)/e_\sigma}$ . Hence we can choose the  $M_n^\sigma$  so that expression (1) equals an arbitrary  $M^\sigma \in \mathbb{C}^{\dim(\sigma) \times \dim(\sigma)}$ .

In particular, for each  $\sigma$ , the  $(i, j)$  entries of the  $M_n^\sigma$  should satisfy

$$\sum_{n \in N} M_n^\sigma[i, j] \begin{pmatrix} \lambda_1(n) \\ \lambda_2(n) \\ \vdots \\ \lambda_{|\mathcal{O}_\ell|}(n) \end{pmatrix} = \begin{pmatrix} M^\sigma[i, j \cdot |\mathcal{O}_\ell|] \\ M^\sigma[i, j \cdot |\mathcal{O}_\ell| + 1] \\ \vdots \\ M^\sigma[i, j \cdot |\mathcal{O}_\ell| + |\mathcal{O}_\ell| - 1] \end{pmatrix} \quad (2)$$

where  $M^\sigma$  on the right-hand-side is blocked into  $d_\sigma \times d_\sigma$  submatrices and indexed accordingly. Thus the values of a given entry of  $M_n^\sigma$  as  $n$  ranges over  $N$ , can be found in an inverse  $N$ -DFT with the appropriate blocks of  $M^\sigma$  as input data.

Observe however that in general,  $\mathcal{O}_\ell$  is a *proper* subset of  $\text{Irr}(H)$ , and hence the aforementioned inverse  $N$ -DFT is underdetermined; for example Equation (2) remains satisfied if we require  $\sum_{n \in N} M_n^\sigma[i, j] \lambda(n) = 0$  for all  $\lambda \in \text{Irr}(H) \setminus \mathcal{O}_\ell$ .

Indeed, we can *simultaneously* solve Equation (2) with respect to several  $\sigma \in \text{Irr}(H)$  via a single inverse  $N$ -DFT, provided the associated orbits  $\mathcal{O}_{i_\sigma}$  are different. To prove the “moreover” part of the theorem statement, then, we set up the following system of equations, for a given  $a \in [r]$ : for each  $\ell$  for which  $f_\ell(\sigma, i, j) = a$  we *simultaneously* require that Equation (2) holds with respect to  $\sigma, i, j$  (and note these are determined by  $a$  since  $f_\ell$  is injective). Since the  $S_\ell$  partition  $\text{Irr}(H)$ , selecting at most one  $\sigma$  from each  $S_\ell$  results in a system that mentions each  $\lambda \in \text{Irr}(N)$  at most once. Hence a single inverse  $N$ -DFT solves this system of equations. See Figure 1. We do this once for each  $a \in [r]$ , to produce the matrices  $M_n^\sigma$  from the original  $M$ , using  $r$  inverse  $N$ -DFTs.  $\square$

### 3.3.1 Choosing the labellings $f_\ell$

To make use of Theorem 3.5, we need to define injective functions  $f_\ell$  from

$$\{(\sigma, i, j) : \sigma \in S_\ell, i \in [\dim(\sigma)/d_\sigma], j \in [e_\sigma]\}$$

to  $[r]$ . We identify the domain of  $f_\ell$  with the entries of a block-diagonal matrix, with rectangular blocks of size  $\dim(\sigma)/d_\sigma \times e_\sigma$ , as  $\sigma$  ranges over  $S_\ell$ . Recall that by Proposition 2.3, the total number of entries in these blocks is  $|H/N|$ .



We will describe functions  $f_\ell$  associating the entries of a block-diagonal matrix of this format (which depends on  $\ell$ ) with a *target* block-diagonal matrix whose format is fixed as follows:

$$\begin{array}{lll}
2 \cdot |H/N| & \text{blocks of size} & 1 \times 1 \\
\lceil 2 \cdot |H/N|/4 \rceil & \text{blocks of size} & 2 \times 2 \\
\lceil 2 \cdot |H/N|/16 \rceil & \text{blocks of size} & 4 \times 4 \\
& \vdots & \\
\lceil 2 \cdot |H/N|/2^{2^i} \rceil & \text{blocks of size} & 2^i \times 2^i \\
& \vdots & \\
2 & \text{blocks of size} & 2^{\lceil \log_2 |H/N| \rceil} \times 2^{\lceil \log_2 |H/N| \rceil}
\end{array}$$

Note that the number of entries of this target matrix is  $O(|H/N| \log |H/N|)$ , and this will be our  $r$ . The association specifying the map  $f_\ell$  is quite simple: we take one column at a time of the source block-diagonal matrix, and if it has height  $w$ , we associate it, top-aligned, with the next-available column among the blocks of size  $2^i \times 2^i$ , for the  $i$  such that  $2^i/2 < w \leq 2^i$ . See Figure 2. Since there can be at most  $|H/N|/w < 2|H/N|/2^i$  columns of height  $w$  in the source matrix (which has  $|H/N|$  entries in total), and the target block-diagonal matrix has at least  $2 \cdot |H/N|/2^i$  columns of width  $i$ , this association is possible.

We will use these mappings when applying Theorem 3.5 to facilitate an efficient “lift” from an intermediate representation to the final  $G$ -DFT. The key benefit of the mappings is that they allow us to combine several matrix-vector products with incompatible formats into one, as illustrated in Figure 2. In order to be able to speak precisely about this combined object, we make the following definition:

**Definition 3.6** (parent matrix). *Given a partition of  $\text{Irr}(H)$  into sets  $S_\ell$ , matrices  $A^\sigma$  with dimensions  $\dim(\sigma)/d_\sigma \times e_\sigma$  (one for each  $\sigma \in \text{Irr}(H)$ ), and functions  $f_\ell$  as above, satisfying*

$$f_\ell(\sigma, i, j) = f_{\ell'}(\sigma', i', j') \Rightarrow A^\sigma[i, j] = A^{\sigma'}[i', j'],$$

define the parent matrix of the  $A^\sigma$  to be the matrix with the format of the target matrix above, and with entry  $(x, y)$  equal to the value of  $A^\sigma[i, j]$  if there exists  $\ell$  for which  $f_\ell(\sigma, i, j) = (x, y)$ , and zero otherwise.

See Figure 3 for an example parent matrix.

### 3.3.2 Computing the intermediate representation

We are at the point now where we can compute the intermediate representation, which we then lift to the final  $G$ -DFT, making critical use of the just-described labellings  $f_\ell$ . The setup is as follows:  $H$  and  $K$  are proper subgroups of group  $G$ , and  $H \cap K = N$  is normal in  $G$ . Let  $X$  be a system of distinct coset representatives of  $N$  in  $H$  and let  $Y$  be a system of distinct coset representatives of  $N$  in  $K$ . Thus  $H = XN$  and  $K = NY$ . Note that  $HK = XNY$  with uniqueness of expression.

When applying the triple subgroup reduction in our final result, it will happen that

$$\frac{|G|}{|HK|} = \frac{|G||N|}{|H||K|}$$

is negligible, and notice that in this case, if  $H$ -DFTs,  $K$ -DFTs, and  $N$ -DFTs have nearly-linear algorithms, then indeed the cost of applying the next lemma is nearly-linear in  $|G|$  as desired.

**Lemma 3.7.** *With  $|Y|$  many  $H$ -DFTs,  $O(|H/N| \log |H/N|) \cdot |Y|$  many inverse  $N$ -DFTs, and  $O(|H/N| \log |H/N|)$  many  $K$ -DFTs, we can compute, from  $\alpha \in \mathbb{C}[G]$  supported on  $HK$ , the following expression:*

$$\sum_{n \in N} \sum_{y \in Y} \bigoplus_{\tau \in \text{Irr}(K)} P_{n,y} \otimes \tau(ny)^T \tag{3}$$

where  $P_{n,y}$  is the parent matrix of the matrices  $\{M_{n,y}^\sigma : \sigma \in \text{Irr}(H)\}$ , and for each  $\sigma, y$ , the  $M_{n,y}^\sigma$  satisfy (with respect to an  $N$ -adapted basis for  $\text{Irr}(H)$ ):

$$\sum_{n \in N} (M_{n,y}^\sigma \otimes J^\sigma) \sigma(n) = \sum_{h \in H} \alpha_{hy} \sigma(h). \quad (4)$$

where  $J^\sigma$  is the  $\dim(\sigma)/e_\sigma \times d_\sigma$  matrix  $(I_{d_\sigma} | I_{d_\sigma} | \cdots | I_{d_\sigma})$  as in Theorem 3.5.

*Proof.* First, compute for each  $y \in Y$  and  $\sigma \in \text{Irr}(H)$  the matrices

$$M_y^\sigma = \sum_{h \in H} \alpha_{hy} \sigma(h),$$

using  $|Y|$  different  $H$ -DFTs. Next, apply Theorem 3.5, once for each  $y$ , to the matrices

$$\bigoplus_{\sigma \in \text{Irr}(H)} M_y^\sigma \in \bigoplus_{\sigma \in \text{Irr}(H)} \mathbb{C}^{\dim(\sigma) \times \dim(\sigma)},$$

together with the labelings  $f_\ell$  from Section 3.3.1, to obtain matrices  $M_{n,y}^\sigma \in \mathbb{C}^{\dim(\sigma)/d_\sigma \times e_\sigma}$  for which

$$\sum_{n \in N} (M_{n,y}^\sigma \otimes J^\sigma) \sigma(n) = M_y^\sigma,$$

at a cost of  $O(|H/N| \log |H/N|) \cdot |Y|$  many inverse  $N$ -DFTs. Note that these  $M_{n,y}^\sigma$  satisfy Equation (4). Let  $P_{n,y}$  be the parent matrix of the matrices  $\{M_{n,y}^\sigma : \sigma \in \text{Irr}(H)\}$ .

For each  $(i, j)$ , the vector  $\beta$  with  $\beta[ny] = P_{n,y}[i, j]$  is an element of  $\mathbb{C}[K]$  and we perform a  $K$ -DFT on it; this entails computing at most  $O(|H/N| \log |H/N|)$  different  $K$ -DFTs because this is the number of entries in the blocks of the block-diagonal matrices  $P_{n,y}$ . At this point we hold, in the aggregate, all of the entries of Expression (3) in the statement of the lemma, and the proof is complete.  $\square$

### 3.3.3 Lifting to a $G$ -DFT

In this section we show how to efficiently lift the intermediate representation, Expression 3 computed via Lemma 3.7, to a  $G$ -DFT. We continue with the notation of the previous section.

Let  $\text{Irr}^*(H)$  denote the *multiset* of irreps of  $H$  that occur in the restrictions of the irreps of  $G$  to  $H$  (with the correct multiplicities), and similarly let  $\text{Irr}^*(K)$  denote the *multiset* of irreps of  $K$  that occur in the restrictions of the irreps of  $G$  to  $K$ . Let  $S$  and  $T$  be the change of basis matrices that satisfy:

$$\begin{aligned} S \left( \bigoplus_{\sigma \in \text{Irr}^*(H)} \sigma(h) \right) S^{-1} &= \bigoplus_{\rho \in \text{Irr}(G)} \rho(h) \quad \forall h \in H \\ T \left( \bigoplus_{\tau \in \text{Irr}^*(K)} \tau(k) \right) T^{-1} &= \bigoplus_{\rho \in \text{Irr}(G)} \rho(k) \quad \forall k \in K. \end{aligned}$$

We further specify that  $S$  should be with respect to an  $N$ -adapted basis for  $\text{Irr}(H)$ .

Notice that for  $n \in N = H \cap K$ , we have:

$$S \left( \bigoplus_{\sigma \in \text{Irr}^*(H)} \sigma(n) \right) S^{-1} = T \left( \bigoplus_{\tau \in \text{Irr}^*(K)} \tau(n) \right) T^{-1},$$

or equivalently

$$\left( \bigoplus_{\sigma \in \text{Irr}^*(H)} \sigma(n) \right) S^{-1} T = S^{-1} T \left( \bigoplus_{\tau \in \text{Irr}^*(K)} \tau(n) \right), \quad (5)$$

a fact we will use shortly.

A  $G$ -DFT with input  $\alpha$  supported on  $HY = HK$  is the expression:

$$\begin{aligned} \sum_{\substack{h \in H \\ y \in Y}} \alpha_{hy} \bigoplus_{\rho \in \text{Irr}(G)} \rho(hy) &= \sum_{y \in Y} \left( \sum_{h \in H} \alpha_{hy} \bigoplus_{\rho \in \text{Irr}(G)} \rho(h) \right) \cdot \left( \bigoplus_{\rho \in \text{Irr}(G)} \rho(y) \right) \\ &= \sum_{y \in Y} S \left( \sum_{h \in H} \alpha_{hy} \bigoplus_{\sigma \in \text{Irr}^*(H)} \sigma(h) \right) S^{-1} T \left( \bigoplus_{\tau \in \text{Irr}^*(K)} \tau(y) \right) T^{-1} \end{aligned}$$

Now for each  $y \in Y$ , the left-most parenthesized expression is an  $H$ -DFT, with certain blocks repeated. By Equation (4) in the statement of Lemma 3.7, each such expression can be rewritten in terms of matrices  $M_{n,y}^\sigma$ , yielding:

$$\begin{aligned} \sum_{\substack{h \in H \\ y \in Y}} \alpha_{hy} \bigoplus_{\rho \in \text{Irr}(G)} \rho(hy) &= \sum_{\substack{y \in Y \\ n \in N}} S \left( \bigoplus_{\sigma \in \text{Irr}^*(H)} (M_{n,y}^\sigma \otimes J^\sigma) \sigma(n) \right) S^{-1} T \left( \bigoplus_{\tau \in \text{Irr}^*(K)} \tau(y) \right) T^{-1} \\ &= \sum_{\substack{y \in Y \\ n \in N}} S \underbrace{\left( \bigoplus_{\sigma \in \text{Irr}^*(H)} (M_{n,y}^\sigma \otimes J^\sigma) \right)}_{(*)} S^{-1} T \left( \bigoplus_{\tau \in \text{Irr}^*(K)} \tau(ny) \right) T^{-1} \quad (6) \end{aligned}$$

where the last line invoked Equation (5) to move  $\sigma(n)$  past  $S^{-1}T$ .

We now focus on Expression (\*). By Proposition 2.1 we can express Expression (\*) as

$$\left( \bigoplus_{\sigma \in \text{Irr}^*(H), \tau \in \text{Irr}^*(K)} ((M_{n,y}^\sigma \otimes J^\sigma) \otimes \tau(ny)^T) \right) \cdot \text{vec}(S^{-1}T) = \text{vec}(*). \quad (7)$$

We next apply two types of simplifications to the block-diagonal matrix on the left. In each, we observe that equalities among blocks allow us to simplify that block-diagonal matrix, at the expense of arranging portions of  $\text{vec}S^{-1}T$  and  $\text{vec}(*)$  into block-diagonal matrices and summing certain entries. The first such observation is that computing

$$\left( \begin{array}{c|c} A & \\ \hline & A \end{array} \right) \cdot \begin{pmatrix} x_1 \\ x_2 \end{pmatrix} = \begin{pmatrix} y_1 \\ y_2 \end{pmatrix}$$

is equivalent to computing  $A \cdot (x_1|x_2) = (y_1|y_2)$ . The second observation is that computing

$$(A|A) \cdot \begin{pmatrix} x_1 \\ x_2 \end{pmatrix} = y$$

is equivalent to computing  $A \cdot (x_1 + x_2) = y$ .

Using the first observation we can thus simplify Equation 7 to:

$$\left( \bigoplus_{\sigma \in \text{Irr}(H), \tau \in \text{Irr}(K)} ((M_{n,y}^\sigma \otimes J^\sigma) \otimes \tau(ny)^T) \right) \cdot X_0 = Y_0,$$

where  $X_0$  is a block-diagonal matrix whose entries coincide with the entries of  $S^{-1}T$ . Next, we notice that  $J^\sigma = I_{d_\sigma} \otimes (1, 1, \dots, 1)$ . The first observation then allows us to simplify Equation 7 futher to:

$$\left( \bigoplus_{\sigma \in \text{Irr}(H), \tau \in \text{Irr}(K)} ((M_{n,y}^\sigma \otimes (1, 1, \dots, 1)) \otimes \tau(ny)) \right) \cdot X_1 = Y_1$$

where again the entries of  $X_1$  coincide with the entries of  $S^{-1}T$ , and the second observation allows us to simplify to:

$$\left( \bigoplus_{\sigma \in \text{Irr}(H), \tau \in \text{Irr}(K)} M_{n,y}^\sigma \otimes \tau(ny)^T \right) \cdot X_2 = Y_2, \quad (8)$$

where now  $X_2$  is a block-diagonal matrix whose entries are sums of entries of  $S^{-1}T$ .

As in the statement of Lemma 3.7, for each  $n, y$ , let  $P_{n,y}$  be the parent matrix of the matrices  $\{M_{n,y}^\sigma : \sigma \in \text{Irr}(H)\}$ . We can rewrite Expression (8) as

$$\left( \bigoplus_{\tau \in \text{Irr}(K)} P_{n,y} \otimes \tau(ny)^T \right) \cdot X_3 = Y_3, \quad (9)$$

where  $X_3$  is a block-diagonal matrix whose entries are sums of entries of  $S^{-1}T$ .

The square blocks of the block-diagonal matrix

$$\left( \bigoplus_{\tau \in \text{Irr}(K)} P_{n,y} \otimes \tau(ny)^T \right)$$

have dimensions  $a_i$  with the property that

$$\sum_i a_i^2 = O(|H/N| \log |H/N|) \cdot |K|,$$

using our earlier accounting for the block sizes of a parent matrix, together with the fact that  $\sum_{\tau \in \text{Irr}(K)} \dim(\tau)^2 = |K|$ . Each  $a_i \times a_i$  block is multiplied by an  $a_i \times w_i$  block of  $X_3$ , to yield an  $a_i \times w_i$  block of the product matrix  $Y_3$ . We now argue that the  $w_i$  satisfy  $\sum_i a_i w_i = 4|G|$ . Each of the two transformations applied to obtain block-diagonal matrices  $Y_0, Y_1$  and then  $Y_2$  preserve the number of entries of the result matrix; these matrices therefore have  $|G|$  entries in the blocks. The final transformation results in a block-diagonal matrix  $Y_3$  which may have *more* entries than  $|G|$ , but this number can be larger by only a factor of four, as illustrated in Figure 2. This is because each column of a block of  $Y_2$  may need to be padded to at most twice its original length, and repeated up to two times (and no more, because the blocks of the  $M_{n,y}^\sigma$  have no more columns than rows, and thus can spill over at most two blocks in the parent matrix). Thus the number of entries in the blocks of  $Y_3$  which equals  $\sum_i a_i w_i$ , is at most  $4|G|$  as stated.

We conclude that the block-matrix multiplication in Expression (9) can be performed efficiently as summarized in the following lemma.

**Lemma 3.8.** *The map from*

$$\sum_{n \in N} \sum_{y \in Y} \bigoplus_{\tau \in \text{Irr}(K)} P_{n,y} \otimes \tau(ny)^T$$

*as computed from input  $\alpha$  supported on  $HY = HK$  in Lemma 3.7, to a  $G$ -DFT with, can be computed at a cost of  $O(|G|^{\omega/2+\epsilon})$  operations, for all  $\epsilon > 0$ .*

*Proof.* We describe how to map a summand  $\bigoplus_{\tau \in \text{Irr}(K)} P_{n,y} \otimes \tau(ny)^T$  to the corresponding summand of Expression (6). This map will be *linear* and will not depend on  $n, y$ , so we apply it once to the entire sum computed by Lemma 3.7, to obtain Expression (6), which is the promised  $G$ -DFT.

We need to perform matrix multiplications of format  $\langle a_i, a_i, w_i \rangle$ , and we know that  $\sum_i a_i^2 = O(|H/N| \log |H/N|) \cdot |K| = L$  and  $\sum_i a_i w_i = 4|G|$ . The cost of such a multiplication is at most

$\max(O(a_i^{\omega+\epsilon}), O(a_i^{\omega-1+\epsilon}w_i))$  for all  $\epsilon > 0$ . Replacing the maximum with a sum, and letting  $a_{\max} = \max_i a_i$ , we obtain an upper bound on the number of operations of

$$\sum_i O(a_i^{\omega+\epsilon}) + O(a_i^{\omega-1+\epsilon}w_i) = O(a_{\max}^{\omega-2+\epsilon}) \sum_i a_i^2 + a_i w_i \leq L^{(\omega-2+\epsilon)/2} \cdot (L + 4|G|). \quad (10)$$

We need to pre-multiply by  $S$  and post-multiply by  $T^{-1}$  to obtain a summand of Expression (6). Both  $S$  and  $T^{-1}$  are block-diagonal with one block for each  $\rho \in \text{Irr}(G)$ , with dimension  $\dim(\rho)$ . Thus the cost of this final pre- and post- multiplication is

$$\sum_{\rho \in \text{Irr}(G)} O(\dim(\rho)^{\omega+\epsilon})$$

which is at most  $O(|G|^{\omega/2+\epsilon})$  by Proposition 2.2 with  $\alpha = \omega + \epsilon$ . The theorem follows from the fact that  $|H||K|/|N| \leq |G|$ , and thus Expression (10) is also upper-bounded by  $O(|G|^{\omega/2+\epsilon})$  (absorbing logarithmic terms into  $|G|^{\epsilon/2}$ ).  $\square$

We now have the main theorem putting together the entire triple subgroup reduction:

**Theorem 3.9** (triple subgroup reduction). *Let  $G$  be a finite group and let  $H, K$  be proper subgroups with  $N = H \cap K$  normal in  $G$ . Then we can compute generalized DFTs with respect to  $G$  at the cost of*

- $|K|/|N|$  many  $H$ -DFTs,
- $O(|H||K|/|N|^2 \log |H/N|)$  many inverse  $N$ -DFTs,
- $O(|H/N| \log |H/N|)$  many  $K$ -DFTs,

plus  $O(|G|^{\omega/2+\epsilon})$  operations, all repeated  $O(|G| \log |G|/|HK|)$  many times, for all  $\epsilon > 0$ .

*Proof.* By Lemma 3.7 we can compute the intermediate representation of a  $G$ -DFT supported on  $HK$ , and applying the map of Lemma 3.8 to this intermediate representation yields a  $G$ -DFT supported on  $HK$ . By Theorem 2.5 we can compute a general  $G$ -DFT at the cost of repeating these two steps  $O(|G| \log |G|/|HK|)$  many times.  $\square$

### 3.3.4 Triple subgroup structure in finite groups

Our main structural theorem on finite groups is the following

**Theorem 3.10.** *There exists a monotone increasing function  $f(x) \leq 2^{c\sqrt{\log x} \log \log x}$  for a universal constant  $c \geq 1$ , such that, for every nontrivial finite group  $G$  one of the following holds*

1.  $G$  has a (possibly trivial) normal subgroup  $N$  and  $G/N$  is cyclic of prime order, or
2.  $G$  has a (possibly trivial) normal subgroup  $N$  and  $G/N$  has proper subgroups  $X, Y$  with  $X \cap Y = \{1\}$  and for which  $|X||N||Y| \geq |G|/f(G)$ .

To connect this theorem to our usage in the previous sections, think of  $H$  as being the subgroup  $\overline{X}N$  and  $K$  as being the subgroup  $N\overline{Y}$ , where  $\overline{X}$  and  $\overline{Y}$  are lifts of  $X$  and  $Y$ , respectively, from  $G/N$  to  $G$ .

*Proof.* Let  $N$  be a maximal normal subgroup of  $G$ . Then  $G/N$  is simple. If it is cyclic of prime order, then we are done. Otherwise we have the following cases, by the Classification Theorem:

1.  $G/N$  is an alternating group  $A_n$  for  $n \geq 5$ . In this case, let  $X$  be the subgroup of  $G/N$  isomorphic to  $A_{n-1}$  and  $Y$  the trivial subgroup of  $G/N$ .

2.  $G/N$  is a finite group of Lie Type. In this case, we refer to Table 4, and we have the following description from Carter [Car89]. For Chevalley and exceptional Chevalley groups, we have that there are subgroups  $B$  and  $U_w^-$  (for each  $w$  in the associated Weyl group  $W$ ) so that elements of  $G/N$  can be expressed *uniquely* as  $bn_wu_w$ , where  $b \in B$ ,  $n_w$  is a lift of  $w \in W$  to  $G$ , and  $u_w \in U_w^-$  (see Corollary 8.4.4 in Carter [Car89]). Uniqueness implies that the conjugate subgroup  $n_wU_w^-n_w^{-1}$  has trivial intersection with  $B$ ; also, by an averaging argument, there exists  $w \in W$  for which  $|Bn_wU_w^-n_w^{-1}| \geq |G/N|/|W|$ . We take  $X = B$  and  $Y = n_wU_w^-n_w^{-1}$ . For twisted Chevalley groups, we have an identical situation (see Corollary 13.5.3 in Carter [Car89]), with subgroup  $B$  replaced by  $B^1$  and subgroup  $U_w^-$  replaced by  $(U_w^-)^1$  (in Carter's notation). Again by an averaging argument there exists  $w \in W$  for which  $|B^1n_w(U_w^-)^1n_w^{-1}| \geq |G/N|/|W|$ , and subgroups  $B^1$  and  $n_w(U_w^-)^1n_w^{-1}$  have trivial intersection; so we take them as our  $X$  and  $Y$ , respectively. Finally we verify from Table 4 that in all cases we have  $f(|G/N|) \geq |W|$ . Thus

$$|X||N||Y| \geq |N||G/N|/|W| \geq |N||G/N|/f(|G/N|) \geq |G|/f(|G|)$$

where we used the fact that  $f$  is increasing.

3.  $G/N$  is a one of the 26 sporadic groups or the Tits group. In this case, we can take  $X = Y = \{1\}$ , by choosing  $c$  in the definition of  $f(x)$  sufficiently large. □

### 3.3.5 Putting it together

Using the structural theorem and the new triple-subgroup reduction recursively, we obtain our final result:

**Theorem 3.11** (main). *For any finite group  $G$ , there is an arithmetic algorithm computing generalized DFTs with respect to  $G$ , using  $O(|G|^{\omega/2+\epsilon})$  operations, for any  $\epsilon > 0$ .*

*Proof.* Fix an arbitrary  $\epsilon > 0$ . Consider the following recursive algorithm to compute a  $G$ -DFT. If  $G$  is trivial then computing a  $G$ -DFT is as well. If  $G$  has a proper subgroup  $H$  of order larger than  $|G|^{1-\epsilon/2}$  then we apply Theorem 3.1 to compute a  $G$ -DFT via several  $H$ -DFTs. Otherwise, applying Theorem 3.10, we obtain a (possibly trivial) normal subgroup  $N$ , and two proper subgroups of  $G$ ,  $H$  and  $K$ , with  $N = H \cap K$ . If  $G/N$  is cyclic of prime order, we apply Corollary 3.3 to compute a  $G$ -DFT via several  $N$ -DFTs. Otherwise, we apply Theorem 3.9 to compute a  $G$ -DFT via several  $H$ -DFTs,  $K$ -DFTs, and inverse  $N$ -DFTs.

Let  $T(n)$  denote an upper bound on the operation count of this recursive algorithm for any group of order  $n$ . We will prove by induction on  $n$ , that there is a universal constant  $C_\epsilon$  for which

$$T(n) \leq C_\epsilon n^{\omega/2+\epsilon} \log n.$$

In the case that we apply Theorem 3.1, the cost is the cost of  $[G : H]$  many  $H$ -DFTs plus  $A_0[G : H]|G|^{\omega/2+\epsilon/2}$  operations (where  $A_0$  is the constant hidden in the big-oh), and by induction this is at most:

$$C_\epsilon [G : H] |H|^{\omega/2+\epsilon} \log |H| + A_0 [G : H] |G|^{\omega/2+\epsilon/2} \leq C_\epsilon |G|^{\omega/2+\epsilon} (\log |G| - 1) + A_0 |G|^{\omega/2+\epsilon}$$

which is indeed less than  $C_\epsilon |G|^{\omega/2+\epsilon} \log |G|$  provided  $C_\epsilon \geq A_0$ .

In the case that we apply Corollary 3.3, our cost is  $p$  many  $N$ -DFTs, plus  $A_1 |G|^{\omega/2+\epsilon}$  operations, which by induction is at most

$$C_\epsilon p (|G|/p)^{\omega/2+\epsilon} \log (|G|/p) + A_1 |G|^{\omega/2+\epsilon} \leq C_\epsilon |G|^{\omega/2+\epsilon} (\log |G| - 1) + A_1 |G|^{\omega/2+\epsilon},$$

which is indeed less than  $C_\epsilon |G|^{\omega/2+\epsilon} \log |G|$  provided  $C_\epsilon \geq A_1$ .

Finally, in the case that we apply Theorem 3.9, let  $A_2$  be the maximum of the constants hidden in the big-ohs in the statement of the Theorem (applied with  $\epsilon/2$ ). Note that by selecting  $C_\epsilon$  sufficiently large, we may assume that  $G$  is sufficiently large, so that two inequalities hold:

$$\begin{aligned} A_2 |H/N| \log |H/N| &\leq \frac{|H/N|^{\omega/2+\epsilon}}{4A_2 f(|G|) \log |G|} \\ |K/N| &\leq \frac{|K/N|^{\omega/2+\epsilon}}{4A_2 f(|G|) \log |G|} \end{aligned}$$

and this is possible because Theorem 3.10 implies that  $|H/N|$  (resp.  $|K/N|$ ) are at least  $|G|^{\epsilon/2}/f(|G|)$ , as otherwise  $|K|$  (resp.  $|H|$ ) would exceed  $|G|^{1-\epsilon/2}$ . Our cost is  $|K/N|$  many  $H$ -DFTs,  $A_2 |H| |K| / |N|^2 \log |H/N|$  many inverse  $N$ -DFTs,  $A_2 |H/N| \log |H/N|$  many  $K$ -DFTs, plus  $A_2 |G|^{\omega/2+\epsilon/2}$  operations, all repeated  $A_2 |G| \log |G| / |HK| \leq A_2 f(|G|) \log |G|$  times. By induction, this is at most

$$\begin{aligned} &\left( C_\epsilon |K/N| |H|^{\omega/2+\epsilon} \log |H| + C_\epsilon A_2 |H| |K| / |N|^2 \log |H/N| |N|^{\omega/2+\epsilon} \log |N| \right. \\ &\left. C_\epsilon A_2 |H/N| \log |H/N| |K|^{\omega/2+\epsilon} \log |K| + A_2 |G|^{\omega/2+\epsilon/2} \right) \cdot A_2 f(|G|) \log |G| \end{aligned}$$

Now the first three summands are each at most

$$\frac{C_\epsilon |G|^{\omega/2+\epsilon} \log |G|}{4A_2 f(|G|) \log |G|}$$

as is the fourth summand provided  $|G|$  is sufficiently large. Thus the entire expression is at most  $C_\epsilon |G|^{\omega/2+\epsilon} \log |G|$ , as required. This completes the proof.  $\square$

## 4 Open problems

Is there a proof of Theorem 3.10 that does not need the Classification Theorem? A second question is whether the dependence on  $\omega$  can be removed. Alternatively, can one show that a running time that depends on  $\omega$  is necessary by showing that an exponent-one DFT for a certain family of groups would imply  $\omega = 2$ ?

**Acknowledgements.** We thank Jonah Blasiak, Tom Church, and Henry Cohn for useful discussions during an AIM SQuaRE meeting.

## References

- [Bau91] Ulrich Baum. Existence and efficient construction of fast Fourier transforms on supersolvable groups. *computational complexity*, 1(3):235–256, Sep 1991.
- [BCS97] P. Bürgisser, M. Clausen, and M. A. Shokrollahi. *Algebraic Complexity Theory*, volume 315 of *Grundlehren der mathematischen Wissenschaften*. Springer-Verlag, 1997.
- [Bet84] Thomas Beth. *Verfahren der schnellen Fourier-Transformation*. Teubner, 1984.
- [Car89] Roger W Carter. *Simple groups of Lie type*, volume 22. John Wiley & Sons, 1989.
- [CB93] Michael Clausen and Ulrich Baum. *Fast Fourier transforms*. Wissenschaftsverlag, 1993.
- [Cla89] Michael Clausen. Fast generalized Fourier transforms. *Theoretical Computer Science*, 67(1):55–63, 1989.

- [CT65] James W. Cooley and John W. Tukey. An algorithm for the machine calculation of complex fourier series. *Mathematics of Computation*, 19(90):297–301, 1965.
- [HJ91] Roger A. Horn and Charles R. Johnson. *Topics in Matrix Analysis*. Cambridge University Press, 1991.
- [HU18a] Chloe Ching-Yun Hsu and Chris Umans. A fast generalized DFT for finite groups of lie type. In Artur Czumaj, editor, *Proceedings of the Twenty-Ninth Annual ACM-SIAM Symposium on Discrete Algorithms, SODA 2018, New Orleans, LA, USA, January 7-10, 2018*, pages 1047–1059. SIAM, 2018.
- [HU18b] Chloe Ching-Yun Hsu and Chris Umans. A new algorithm for fast generalized DFTs. *CoRR*, abs/1707.00349v3, 2018. Full version of [HU18a].
- [Lev92] Arieh Lev. On large subgroups of finite groups. *Journal of Algebra*, 152(2):434–438, 1992.
- [Mas98] David Keith Maslen. The efficient computation of Fourier transforms on the symmetric group. *Math. Comput.*, 67(223):1121–1147, 1998.
- [MR97a] David Maslen and Daniel Rockmore. Separation of variables and the computation of Fourier transforms on finite groups, I. *Journal of the American Mathematical Society*, 10(1):169–214, 1997.
- [MR97b] David K Maslen and Daniel N Rockmore. Generalized FFTs – a survey of some recent results. In *Groups and Computation II*, volume 28, pages 183–287. American Mathematical Soc., 1997.
- [MRW16] David Maslen, Daniel N Rockmore, and Sarah Wolff. The efficient computation of Fourier transforms on semisimple algebras. *arXiv preprint arXiv:1609.02634*, 2016. To appear in *Journal of Fourier Analysis and Applications*.
- [Roc95] Daniel N. Rockmore. Fast Fourier transforms for wreath products. *Applied and Computational Harmonic Analysis*, 2(3):279 – 292, 1995.
- [Roc02] Daniel N Rockmore. Recent progress and applications in group FFTs. In *Signals, Systems and Computers, 2002. Conference Record of the Thirty-Sixth Asilomar Conference on*, volume 1, pages 773–777. IEEE, 2002.
- [Wik17] Wikipedia. List of finite simple groups — wikipedia, the free encyclopedia, 2017. [Online; accessed 30-June-2017].



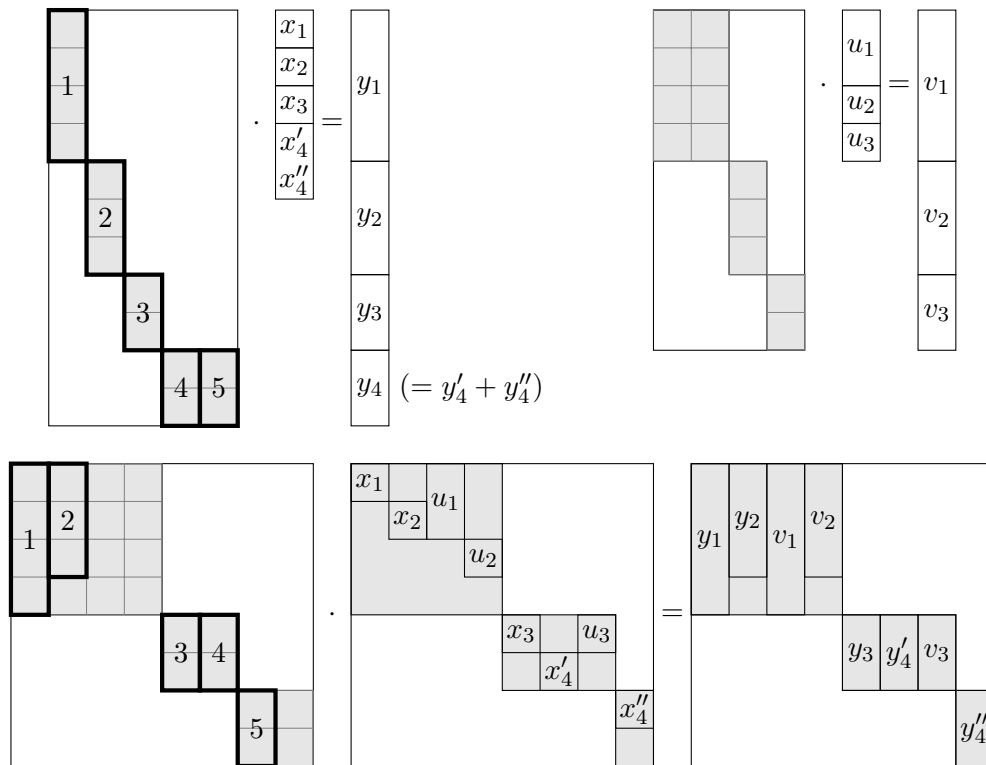


Figure 2: How the  $f_\ell$  functions are defined and used. The bold columns of the block-diagonal matrix in the upper-left are associated to the columns of the target block-diagonal matrix on the bottom-left. The columns of the block-diagonal matrix in the upper-right are also associated the manner described in Section 3.3.1, although this association is not shown in the figure. We see that the two matrix-vector multiplications at the top can be combined into the single matrix product on the bottom, provided that similarly labeled entries of the two source matrices are guaranteed to contain identical values. Unlabeled cells of the middle-bottom matrix contain zeros. Note that in the bottom-right matrix each segment of the original vectors  $y$  and  $v$  may be padded up to twice its original length (but not more), and it may be repeated up to twice and summed (as  $y'_4$  and  $y''_4$  are) if the columns of the associated block are mapped to two different blocks in the target matrix. More than two repetitions are not possible because the source blocks all have at most as many columns as rows.

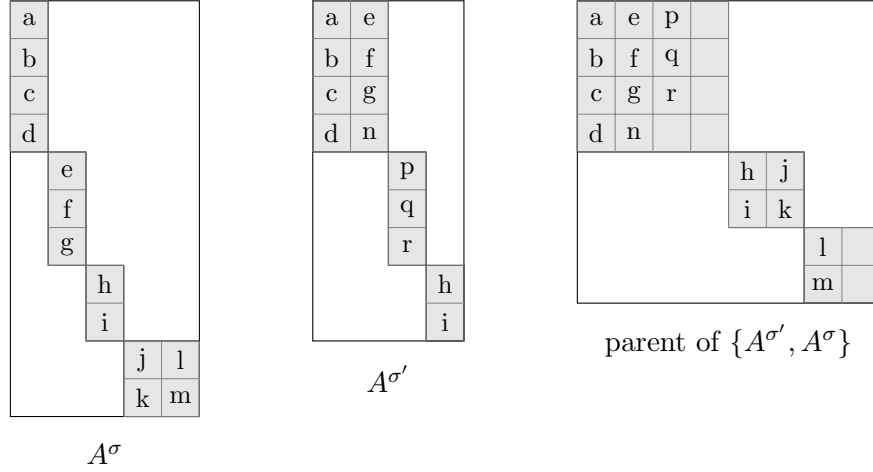


Figure 3: An example parent matrix. Unlabeled entries are zero.

Name	Family	$ W $	$ G $
Chevalley	$A_\ell(q)$	$(\ell + 1)!$	$q^{\Theta(\ell^2)}$
	$B_\ell(q)$	$2^\ell \ell!$	$q^{\Theta(\ell^2)}$
	$C_\ell(q)$	$2^\ell \ell!$	$q^{\Theta(\ell^2)}$
	$D_\ell(q)$	$2^{\ell-1} \ell!$	$q^{\Theta(\ell^2)}$
Exceptional Chevalley	$E_6(q)$	$O(1)$	$q^{\Theta(1)}$
	$E_7(q)$	$O(1)$	$q^{\Theta(1)}$
	$E_8(q)$	$O(1)$	$q^{\Theta(1)}$
	$F_4(q)$	$O(1)$	$q^{\Theta(1)}$
	$G_2(q)$	$O(1)$	$q^{\Theta(1)}$
Steinberg	${}^2A_\ell(q^2)$	$2^{\lceil \ell/2 \rceil} \lceil \ell/2 \rceil!$	$q^{\Theta(\ell^2)}$
	${}^2D_\ell(q^2)$	$2^{\ell-1}(\ell-1)!$	$q^{\Theta(\ell^2)}$
	${}^2E_6(q^2)$	$O(1)$	$q^{\Theta(1)}$
	${}^3D_4(q^3)$	$O(1)$	$q^{\Theta(1)}$
Suzuki	${}^2B_2(q), q = 2^{2n+1}$	$O(1)$	$q^{\Theta(1)}$
Ree	${}^2F_4(q), q = 3^{2n+1}$	$O(1)$	$q^{\Theta(1)}$
	${}^2G_2(q), q = 3^{2n+1}$	$O(1)$	$q^{\Theta(1)}$

Figure 4: Families of finite groups  $G$  of Lie type, together with the size of their associated Weyl group  $W$ . These include all simple finite groups other than cyclic groups, the alternating groups, the 26 sporadic groups, and the Tits group. See [Lev92, Wik17, Car89] for sources. The Suzuki, Steinberg and Ree families are also called the *twisted Chevalley* groups.