

Improved bound for complexity of matrix multiplication

A. M. Davie and **A. J. Stothers**

School of Mathematics, University of Edinburgh, King's Buildings,
Mayfield Road, Edinburgh EH9 3JZ, UK (a.davie@ed.ac.uk)

(MS received 14 November 2011; accepted 29 February 2012)

We give a new bound $\omega < 2.37369$ for the exponent of complexity of matrix multiplication, giving a small improvement on the previous bound obtained by Coppersmith and Winograd. The proof involves an extension of the method used by these authors. We have attempted to make the exposition self-contained.

1. Introduction

In 1968 Strassen [16] described a method of multiplying two 2×2 matrices using only seven multiplications. By applying it recursively it was shown that two $n \times n$ matrices could be multiplied using $O(n^\rho)$ operations in all, where $\rho = \log_2 7 \approx 2.81$. This discovery led to a period of activity devoted to the improvement of this bound. One can define an exponent of complexity ω , which is, in effect, the smallest number such that two $n \times n$ matrices can be multiplied using $O(n^{\omega+\epsilon})$ operations using a method of the type described by Strassen (a precise definition is given in § 2). Then Strassen's result says that $\omega \leq \log_2 7$, so it is natural to wish to determine the true value of ω . In the decade or so after 1978, several successive reductions in the upper bound for ω were achieved [2, 3, 6, 8–10, 12, 14, 17], leading to the bound $\omega < 2.375477$ obtained by Coppersmith and Winograd [7]. It is an open question whether $\omega = 2$.

We present a small improvement of the Coppersmith–Winograd bound, namely, $\omega < 2.3736897$, which was first obtained in Stothers [15]. The proof uses an extension of the ideas of [7] and relies heavily on a combinatorial construction from [7], which depends on probabilistic arguments and a theorem of Salem and Spencer on sets of integers containing no three terms in arithmetical progression.

We have tried to organize the proof so that it is reasonably concise but also self-contained, to make it more accessible, in the belief that at least parts of the argument may be of wider interest outside the algebraic complexity community. To this end we derive the results we need from algebraic complexity theory in § 2, which also contains general background material. In § 3 we present the essence of the combinatorial construction as a lemma, which we subsequently apply several times. We first apply it in § 4 to recover the Coppersmith–Winograd bound. Although our presentation is rather different, the proof is essentially the same as that of [7], except for the proof of lemma 4.1, for which we give a somewhat simpler argument. Using § 4 as preparation and motivation, we obtain our new bound in § 5. The argument is similar but more complicated and technical.

Much more material on this subject and information on the history can be found in Burgisser *et al.* [5] and a survey by Pan [11].

2. Background

We fix a field \mathbb{F} and consider a trilinear form ϕ defined on $U \times V \times W$, where U , V and W are finite-dimensional vector spaces over \mathbb{F} . Note that such a ϕ gives, in a natural way, a bilinear mapping from $U \times V \rightarrow W^*$, where W^* is the dual of W , and so we have a natural one-to-one correspondence between such trilinear forms and bilinear mappings $U \times V \rightarrow W^*$.

We are interested in the case $U = V = W = M_n$, where $M_n = M_n(\mathbb{F})$ is the space of $n \times n$ matrices with entries in \mathbb{F} . Using the trace, we can identify M_n^* with M_n and then matrix multiplication $(A, B) \rightarrow AB$ is a bilinear mapping from $M_n \times M_n \rightarrow M_n$ and the corresponding trilinear form is \mathcal{M}_n given by $\mathcal{M}_n(A, B, C) = \text{tr}(ABC)$. Occasionally, we consider non-square matrix multiplications, and refer to the multiplication of an $m \times n$ matrix with an $n \times p$ matrix as an (m, n, p) matrix product.

The rank $R(\phi)$ of a trilinear form ϕ is defined as the smallest value of $r \in \mathbb{N}$ such that there exist $\rho_q \in U^*$, $\sigma_q \in V^*$ and $\theta_q \in W^*$ for $q = 1, \dots, r$ such that

$$\phi(u, v, w) = \sum_{q=1}^r \rho_q(u) \sigma_q(v) \theta_q(w)$$

for all $u \in U$, $v \in V$ and $w \in W$. We write $r(n)$ for $R(\mathcal{M}_n)$, the rank of multiplication of $n \times n$ matrices. Then $r(n)$ is the smallest r such that there exist ρ_{qij} , σ_{qij} , θ_{qij} in \mathbb{F} for $q = 1, \dots, r$ and $i, j = 1, \dots, n$ such that

$$\sum_{q=1}^r \rho_{qij} \sigma_{qkl} \theta_{qst} = \begin{cases} 1 & \text{if } j = k, l = s, t = i, \\ 0 & \text{otherwise} \end{cases} \quad (2.1)$$

for $i, j, k, l, s, t = 1, \dots, n$.

Given trilinear forms ϕ on $U \times V \times W$ and ψ on $U' \times V' \times W'$, we can define the tensor product $\phi \otimes \psi$ on $(U \otimes U') \times (V \otimes V') \times (W \otimes W')$, and it is elementary that $R(\phi \otimes \psi) \leq R(\phi)R(\psi)$. Applied to \mathcal{M}_m and \mathcal{M}_n , this gives $r(mn) \leq r(m)r(n)$ for $m, n \in \mathbb{N}$. From this, and the elementary fact that $r(n)$ is a non-decreasing function of n , we deduce that $\log r(n)/\log n$ converges to a limit ω as $n \rightarrow \infty$, and that

$$r(n) \geq n^\omega \quad \text{for all } n. \quad (2.2)$$

We now relate ω to the number of arithmetical operations (addition or multiplication of two elements of \mathbb{F}) needed to multiply two matrices in M_n . To do this, consider two matrices $A, B \in M_{km}$, and write A in block form (A_{ij}) , $i, j = 1, \dots, m$ where $A_{ij} \in M_k$, and B similarly. Then, using (2.1), we get

$$C_{ts} = \sum_{q=1}^{r(m)} \theta_{qst} \left(\sum_{i,j=1}^m \rho_{qij} A_{ij} \right) \left(\sum_{k,l=1}^m \sigma_{qkl} B_{kl} \right)$$

for the block form of the product $C = AB$. From this we see that, if two $k \times k$ matrices can be multiplied with $P(k)$ operations, then $mk \times mk$ matrices can be

multiplied with $P(mk)$ operations, where $P(mk) \leq r(m)P(k) + K(m)k^2$ so, by induction, $P(m^j) \leq K'(m)r(m)^j$, where $K(m)$, etc., denote constants depending only on m . Hence, $P(n) \leq K''(m)n^\alpha$ where $\alpha = \log r(m)/\log m$.

It follows that, for any $\epsilon > 0$, there exist C such that, for any $n \in \mathbb{N}$, two $n \times n$ matrices can be multiplied in $Cn^{\omega+\epsilon}$ operations. This motivates the quest for bounds for ω that started with Strassen's proof [16] that $r(2) \leq 7$, which implies that $\omega \leq \log_2 7$ and that $n \times n$ matrices can be multiplied in $O(n^{\log_2 7})$ operations. It is trivial that $\omega \geq 2$ and it is still an open problem whether $\omega = 2$.

Several successive improvements to Strassen's upper bound have been found. While Strassen's bound is based on a bound for $r(2)$, it has proved difficult to get good estimates for $r(n)$ for other small values of n , and the subsequent improvements have been based instead on the development of methods for getting good asymptotic bounds for $r(n)$ for large n . We now describe the results that we require from this theory.

First, some notation. If ϕ is a trilinear form and $N \in \mathbb{N}$, then $\oplus^N \phi$ denotes the direct sum of N copies of ϕ , and ϕ^N denotes the tensor product of N copies of ϕ .

We start with a technical lemma.

LEMMA 2.1. *For any $n, k \in \mathbb{N}$, we have $r(nk) \leq R(\oplus^{r(n)} \mathcal{M}_k)$.*

Proof. Write $r = r(n)$ and $R = R(\oplus^r \mathcal{M}_k)$. The idea is that we can reduce the multiplication of two matrices in M_{nk} to r multiplications of pairs of matrices in M_k , and this we can do with R multiplications (in \mathbb{F}).

To make this precise, from the definition of R , we can find linear functionals $\phi_{iq}, \psi_{iq}, \theta_{iq}$ on M_k for $i = 1, \dots, r$ and $q = 1, \dots, R$ such that if $A_i, B_i, C_i \in M_k$, then

$$\sum_{i=1}^r \text{tr}(A_i B_i C_i) = \sum_{q=1}^R \left(\sum_i \phi_{iq}(A_i) \right) \left(\sum_i \psi_{iq}(B_i) \right) \left(\sum_i \theta_{iq}(C_i) \right) \quad (2.3)$$

Now let ρ_{qjk} , etc., be as given by (2.1) and let $A, B, C \in M_{nk}$. Partition A as $(A_{jl})_{j,l=1}^n$, where $A_{jl} \in M_k$, and define

$$\phi_q(A) = \sum \rho_{ijk} \phi_{iq}(A_{jk}),$$

and similarly for B and C . By applying (2.3) to $A_i = \sum \alpha_{ijk} A_{jk}$, etc., we see that

$$\text{tr}(ABC) = \sum_{q=1}^R \phi_q(A) \psi_q(B) \theta_q(C)$$

and the result follows. □

A crucial tool is the notion of *border rank*, introduced in [3], which we now define. This is a modified rank, using trilinear forms with values in the ring $\mathcal{F}[\lambda]$ of polynomials over \mathcal{F} in an indeterminate λ , which is generally smaller than the rank defined above and is better for obtaining asymptotic bounds.

Let U, V and W be vector spaces over \mathbb{F} and let $\phi: U \times V \times W \rightarrow \mathbb{F}$ be a trilinear form. If $q, r \in \mathbb{N}$, we write $\phi \leq_q r$ if there is a trilinear mapping $\psi: U \times V \times W \rightarrow \mathbb{F}[\lambda]$

and, for $i = 1, \dots, r$, linear mappings u_i, v_i and w_i from U, V and W , respectively, to $\mathbb{F}[\lambda]$ such that

$$\lambda^{q-1}\phi(x, y, z) + \lambda^q\psi(x, y, z) = \sum_{i=1}^r u_i(x)v_i(y)w_i(z). \quad (2.4)$$

Then we define the border rank $\bar{R}(\phi)$ as the smallest r such that, for some $q \in \mathbb{N}$, we have $\phi \trianglelefteq_q r$.

Trivially, we have $\bar{R}(\phi) \leq R(\phi)$, and the inequality is often strict. In the opposite direction, if $\phi \trianglelefteq_q r$, then, by expanding u_i, v_i and w_i in powers of λ , one can express the coefficient of λ^{q-1} in $u_i(x)v_i(y)w_i(z)$ as a sum of $\frac{1}{2}q(q+1)$ products, and so $R(\phi) \leq \frac{1}{2}rq(q+1) \leq rq^2$. This bound is most effective when applied to a large tensor power, as we now show.

LEMMA 2.2. *If $\phi \trianglelefteq_q r$, then $R(\phi^N) \leq N^2q^2r^N$ for any $N \in \mathbb{N}$.*

Proof. By taking the N th tensor power of (2.4), we obtain a similar equation to (2.4) for $\lambda^{N(q-1)}\phi^N$, with r replaced by r^N on the right. Hence, $\phi^N \trianglelefteq_{N(q-1)+1} r^N$, and from the previous paragraph, $R(\phi^N) \leq N^2q^2r^N$. \square

One can think of the border rank as an estimate of the ease of computation of ϕ^N . To apply this to matrix multiplication, we next define a complementary notion of *value* of ϕ , which measures how useful ϕ^N is for computing matrix products.

If ϕ is a trilinear form, then, for $\rho \in [2, 3]$ and $N \in \mathbb{N}$, we define $V_{\phi, N}(\rho) = \sup(kn^\rho)^{1/3N}$, where the supremum is over all choices of $k, n \in \mathbb{N}$ such that $\oplus^k \mathcal{M}_n$ is isomorphic to a restriction of $(\phi \otimes \pi\phi \otimes \pi^2\phi)^N$, where π denotes the action of a cyclic permutation of (x, y, z) . Then it is straightforward that $V_{\phi, Nr}(\rho) \geq V_{\phi, N}(\rho)$ for $r \in \mathbb{N}$ and also that $V_{\phi, N}(\rho)^N$ is increasing with N . It then follows that $V_{\phi, N}(\rho)$ converges to a limit $V_\phi(\rho)$ as $N \rightarrow \infty$, and that $V_\phi(\rho) \geq V_{\phi, N}(\rho)$ for all N . The function V_ϕ is the *value* of ϕ . This definition of value is a modification of that given in [7] which is more suitable for our approach; one can show that, in fact, it is equivalent to the definition in [7].

One should think of ρ as a candidate value for ω (this is the reason for the range $[2, 3]$). The basic result we shall use is as follows.

PROPOSITION 2.3. *For any trilinear form ϕ we have $V_\phi(\omega) \leq \bar{R}(\phi)$.*

Proof. Write $r = \bar{R}(\phi)$ and $V = V_\phi(\omega)$. Then $\phi \trianglelefteq_q r$ for some $q \in \mathbb{N}$. Let $\epsilon > 0$. Then for any sufficiently large N we can find $k, n \in \mathbb{N}$ so that $kn^\omega \geq V^{3N(1-\epsilon)}$ and $\oplus^k \mathcal{M}_n$ is isomorphic to a restriction of $(\phi \otimes \pi\phi \otimes \pi^2\phi)^N$, so that, by lemma 2.2, we have $R(\oplus^k \mathcal{M}_n) \leq 9q^2N^2r^{3N}$.

Now there is $C > 0$ such that $r(m) \leq Cm^{\omega+\epsilon}$ for all $m \in \mathbb{N}$. Then we can find $m \in \mathbb{N}$ so that $r(m) \leq k$ and $Cm^{\omega+\epsilon} \geq \frac{1}{2}k$. Then, by (2.2) and lemma 2.1, we have

$$(mn)^\omega \leq r(mn) \leq R(\oplus^k \mathcal{M}_n) \leq 9q^2N^2r^{3N}.$$

Then

$$V^{3N(1-\epsilon)} \leq kn^\omega \leq 2C(mn)^{\omega+\epsilon} \leq 2C(9q^2N^2r^{3N})^{1+\epsilon/\omega}.$$

Taking $(3N)$ th roots, and letting $N \rightarrow \infty$, we get $V^{1-\epsilon} \leq r^{1+\epsilon/\omega}$ for all $\epsilon > 0$, so $V \leq r$ as required. \square

Proposition 2.3 is a special case of Schönhage’s asymptotic sum inequality [14] which allows direct sums of matrix products of different shapes and sizes.

To apply proposition 2.3 we use a trilinear form χ from [7], defined as follows. Fix a positive integer $q > 1$, and label vectors $x \in \mathbb{F}^{q+2}$ by $x = (x^{[0]}, x_1^{[1]}, \dots, x_q^{[1]}, x^{[2]})$. For $x, y, z \in \mathbb{F}^{q+2}$ we define

$$\begin{aligned} \phi_{011}(x, y, z) &= x^{[0]} \sum_{i=1}^q y_i^{[1]} z_i^{[1]}, \\ \phi_{101}(x, y, z) &= y^{[0]} \sum_{i=1}^q x_i^{[1]} z_i^{[1]}, \\ \phi_{110}(x, y, z) &= z^{[0]} \sum_{i=1}^q y_i^{[1]} x_i^{[1]}, \\ \phi_{200}(x, y, z) &= x^{[2]} y^{[0]} z^{[0]}, \\ \phi_{020}(x, y, z) &= y^{[2]} x^{[0]} z^{[0]}, \\ \phi_{002}(x, y, z) &= z^{[2]} y^{[0]} x^{[0]}. \end{aligned}$$

Then

$$\chi = \phi_{011} + \phi_{101} + \phi_{110} + \phi_{200} + \phi_{020} + \phi_{002} \tag{2.5}$$

is a trilinear form on \mathbb{F}^{q+2} .

We have

$$\begin{aligned} \lambda^3 \chi(x, y, z) &= (1 - q\lambda) \sigma(x) \sigma(y) \sigma(z) - \tau(x) \tau(y) \tau(z) \\ &\quad + \lambda \sum_{i=1}^q \psi_i(x) \psi_i(y) \psi_i(z) + \lambda^4 \mu(x, y, z), \end{aligned}$$

where $\sigma(x) = x^{[0]} + \lambda^3 x^{[2]}$, $\tau(x) = x^{[0]} + \lambda^2 \sum_{i=1}^q x_i^{[1]}$, $\psi_i(x) = x^{[0]} + \lambda x_i^{[1]}$ and μ is a trilinear mapping with values in $\mathbb{F}[\lambda]$. It follows that $\bar{R}(\chi) \leq q + 2$.

Then, if we have a lower bound for $V_\chi(\rho)$, we can use proposition 2.3 to obtain an upper bound for ω . More precisely, if we can show that $V_\chi(\rho) \geq f(\rho)$ for some strictly increasing continuous function f on $[2, 3]$, then we deduce that $\omega \leq \rho_0$, where ρ_0 is the unique solution of $f(\rho_0) = q + 2$.

In the rest of the paper we obtain lower bounds for V_χ . The idea is to take a large tensor power of χ and then, by setting certain blocks of variables to zero, to obtain a restriction that is a direct sum of matrix products. This is really a matter of combinatorics, and the required combinatorial lemma is treated in the next section, followed by the value bounds and consequent bounds of ω in the last two sections.

We note that, in the opposite direction, finding lower bounds for the rank $r(n)$ seems to be difficult (see [4]).

3. Combinatorial lemma

In this section we prove a lemma which is essentially a distillation of the combinatorial argument introduced and used repeatedly by Coppersmith and Winograd [7].

If $M > 1$ is an integer, then, as usual, we denote the rings of integers modulo M by \mathbb{Z}_M . \mathbb{P} denotes probability. We start with an elementary observation.

LEMMA 3.1. *Suppose A is an $m \times r$ matrix with entries in \mathbb{Z}_M such that the mapping $u \rightarrow Au$ maps \mathbb{Z}_M^r onto \mathbb{Z}_M^m . Let w be a random vector in \mathbb{Z}_M^r with $\mathbb{P}(w = v) = M^{-r}$ for each $v \in \mathbb{Z}_M^r$. Then $\mathbb{P}(Aw = u) = M^{-m}$ for every $u \in \mathbb{Z}_M^m$.*

In other words, the components of Aw are independent random variables, each having a uniform distribution on \mathbb{Z}_M .

Proof. Given $u \in \mathbb{Z}_M^m$, we can find $v \in \mathbb{Z}_M^r$ with $Av = u$ and then $\mathbb{P}(Aw = u) = \mathbb{P}(A(w - v) = 0) = \mathbb{P}(Aw = 0)$, since the random vectors w and $w - v$ have the same distribution. So $\mathbb{P}(Aw = u)$ is independent of u and the result follows. \square

We will need the Salem–Spencer theorem [13], which we now state. A set B of integers is called a *Salem–Spencer set* if, whenever $a, b, c \in B$ with $a + b = 2c$, we must have $a = b = c$. Then the Salem–Spencer result is as follows.

PROPOSITION 3.2. *Given $\epsilon > 0$, there is a positive constant C such that if $M \in \mathbb{N}$, one can find a Salem–Spencer set $B \subseteq \{1, 2, \dots, M\}$ with $|B| \geq CM^{1-\epsilon}$.*

A proof can be found in [1].

We now require some notation. For an integer $d \geq 2$, let \mathbb{H}_d be the set $\{0, 1, \dots, d\}$ and let Ω_d be the set of triples $(a, b, c) \in \mathbb{H}_d^3$ such that $a + b + c = d$. For $n \in \mathbb{N}$, we can form the n -fold Cartesian products \mathbb{H}_d^n and Ω_d^n . Then Ω_d^n can be regarded as the set of triples $(I, J, K) \in (\mathbb{H}_d^n)^3$ such that $I_j + J_j + K_j = d$ for each $j = 1, 2, \dots, n$. For $i = 1, 2, 3$, we define $P_i: \Omega_d^n \rightarrow \mathbb{H}_d^n$ by $P_1(I, J, K) = I$, $P_2(I, J, K) = J$ and $P_3(I, J, K) = K$. Now we can state the main result of this section.

LEMMA 3.3. *Given an integer $d \geq 2$ and $\epsilon > 0$, we can find $C > 0$ such that the following holds. Suppose that we have $n, R \in \mathbb{N}$ and sets $S_0 \subseteq S \subseteq \Omega_d^n$ and, for $i = 1, 2, 3$, $F_i \subseteq \mathbb{H}_d^n$ satisfying $P_i(S) = F_i$ and $|P_i^{-1}(I) \cap S| \leq R$ for all $I \in F_i$. Then there are subsets G_i of F_i for $i = 1, 2, 3$ such that, if we write $T = (G_1 \times G_2 \times G_3) \cap S$, then $T \subseteq S_0$, P_i maps T one-to-one onto G_i and $|T| \geq CR^{-1-\epsilon}|S_0|$.*

Proof. Let K be the product of all prime numbers less than or equal to d . We can find an integer M with $4R < M \leq 4R + K$ such that M is prime to K (and, hence, to any number in $\{2, \dots, d\}$). Then we can find a Salem–Spencer set B in $\{1, 2, \dots, 2R\}$ with $|B| \geq C_1 R^{1-\epsilon}$. We think of B as a subset of \mathbb{Z}_M and note that, as $2R < \frac{1}{2}M$, if $a, b, c \in B$ with $a + b = 2c \pmod{M}$, then $a = b = c$.

Now we consider independent random variables $w_0, w_*, w_1, w_2, \dots, w_n$ taking values in \mathbb{Z}_M , each with uniform distribution. Then, for $I \in \mathbb{H}_d^n$, we define \mathbb{Z}_M -valued random variables $w^i(I)$ for $i = 1, 2, 3$ as follows ($\frac{1}{2}$ being well-defined in \mathbb{Z}_M as M is odd):

$$w^1(I) = w_0 + \sum_{j=1}^n w_j I_j,$$

$$w^2(I) = w_* + \sum_{j=1}^n w_j I_j,$$

$$w^3(I) = \frac{1}{2} \left(w_0 + w_* + \sum_{j=1}^n (d - I_j) w_j \right).$$

Note that if $(I, J, K) \in \Omega_d^n$, then $w^1(I) + w^2(J) = 2w^3(K)$. Note that, from lemma 3.3, if $I, J, J' \in \mathbb{H}^n$ and $J \neq J'$, then $w^1(I), w^2(J)$ and $w^2(J')$ are mutually independent, each being uniformly distributed on \mathbb{Z}_M .

Then, for $i = 1, 2, 3$, we define

$$H_i = \{I \in F_i : w^i(I) \in B\}, \quad R_i = \{I \in H_i : |P_i^{-1}(I) \cap S| > 1\}, \quad H_i^0 = H_i \setminus R_i.$$

Then, for fixed $\psi = (I, J, K) \in S$, we have that $\psi \in H_1 \times H_2 \times H_3$ if and only if $w^1(I), w^2(J)$ and $w^3(K)$ are all in B , which, in view of the Salem–Spencer property, is equivalent to the existence of $b \in B$ such that $w^1(I) = w^2(J) = b$ (and then automatically $w^3(K) = b$). For each $b \in B$, the probability that $w^1(I) = w^2(J) = b$ is M^{-2} . Hence,

$$\mathbb{P}(\psi \in H_1 \times H_2 \times H_3) = |B|M^{-2}.$$

Similarly, for $\psi \in S$, the statement that $\psi \in H_1 \times H_2 \times H_3$ and $I \in R_1$ is equivalent to the existence of $b \in B$ and J' and K' in \mathbb{H}_d^n such that $w^1(I) = w^2(J) = w^2(J') = b$, $J' \neq J$ and $(I, J', K') \in S$. Now, for $\psi \in S$, and for each of the $\leq R - 1$ choices of J', K' such that $J' \neq J$ and $(I, J', K') \in S$, we have $\mathbb{P}(w^1(I) = w^2(J) = w^2(J') = b) = M^{-3}$. The same applies for $i = 2, 3$, and we obtain for each i that

$$\mathbb{P}(\psi \in H_1 \times H_2 \times H_3 \quad \text{and} \quad P_i(\psi) \in R_i) \leq |B|RM^{-3} \leq \frac{1}{4}|B|M^{-2}.$$

Hence,

$$\mathbb{P}(\psi \in H_1^0 \times H_2^0 \times H_3^0) \geq |B|M^{-2} - 3 \times \frac{1}{4}|B|M^{-2} = \frac{1}{4}|B|M^{-2}.$$

Now, if we define $T = S_0 \cap (H_1^0 \times H_2^0 \times H_3^0)$, then the expectation $\mathbb{E}|T| \geq \frac{1}{4}|B|M^{-2}|S_0|$. So there is a choice of $w_0, w_*, w_1, \dots, w_n$ such that

$$|T| \geq \frac{1}{4}|B|M^{-2}|S_0| \geq \frac{1}{4}CR^{1-\epsilon}(8R)^{-2}|S_0| = C'R^{-1-\epsilon}|S_0|.$$

Finally, set $G_i = H_i^0 \cap P_i^{-1}(T)$, and the stated properties hold. \square

We make a few remarks on our application of this lemma, which we will use several times in §§ 4 and 5. d always takes one of the values 2, 4 or 8. And we will always have $|P_i^{-1} \cap S| = R$ for $I \in F_i$, which implies $|F_i| = |S|/R$ for each i , and also $R \leq |F_1|$. Then the inequality in the conclusion of the lemma gives

$$|T| \geq C|F_1|^{1-\epsilon}|S_0|/|S|, \tag{3.1}$$

which will be more convenient to use. In the simpler applications we in fact have $S_0 = S$, but in others, the estimation of $|S_0|/|S|$ will be a significant part of the argument.

The choices of S and S_0 needed for our applications fall into two types, which we now discuss further.

3.1. Type 1

Let $W_d^{(n)}$ denote the set of $\Delta = (\Delta_\mu)_{\mu \in \Omega_d}$ in which the Δ_μ are non-negative integers such that $\sum_{\mu \in \Omega_d} \Delta_\mu = n$. For $i = 1, 2, 3$, we define $Q_i \Delta = (u_0, \dots, u_d)$, where $u_k = \sum \Delta_\mu$, the sum being over those $\mu \in \Omega_d$ such that $\mu_i = k$. For example, in the case $d = 2$, we have

$$Q_2 \Delta = (\Delta_{101} + \Delta_{002} + \Delta_{200}, \Delta_{011} + \Delta_{110}, \Delta_{020}).$$

We denote by $S(\Delta)$ the set of elements of Ω_d^n having Δ_μ occurrences of μ for each $\mu \in \Omega_d$.

Now fix $\Gamma \in W_d^{(n)}$ and define

$$S_0 = S(\Gamma), \quad F_i = P_i(S_0), \quad S = \bigcap_{i=1}^3 P_i^{-1}(F_i) = \bigcup_{\Delta \in \Lambda_\Gamma} S(\beta),$$

where Λ_Γ is the set of all $\Delta \in W_d^{(n)}$ such that $Q_i \Delta = Q_i \Gamma$ for $i = 1, 2, 3$. We assume the symmetry of Γ , that it is invariant under permutations, so that, for example, $\Gamma_{103} = \Gamma_{310}$, etc. This implies that $F_1 = F_2 = F_3$, and that $Q_1 \Gamma = Q_2 \Gamma = Q_3 \Gamma$. We let $(A_0, A_1, \dots, A_d) = Q_1 \Gamma$. Note that Λ_Γ may contain non-symmetric β .

Now we have

$$|S_0| = \frac{n!}{\prod_{\mu \in \Omega_d} \Gamma_\mu!} \quad \text{and} \quad |S| = \sum_{\Delta \in \Lambda_\Gamma} \frac{n!}{\prod_{\mu \in \Omega_d} \Delta_\mu!}.$$

The number of terms in this sum is less than or equal to $n^{m(d)}$, where $m(d) = |\Omega_d|$, and hence

$$|S|/|S_0| = \sum_{\Delta \in \Lambda_\Gamma} \prod_{\mu \in \Omega_d} (\Delta_\mu!/\Gamma_\mu!) \leq n^{m(d)} \max_{\Delta \in \Lambda_\Gamma} \prod_{\mu \in \Omega_d} (\Delta_\mu!/\Gamma_\mu!). \tag{3.2}$$

Note also that F_1 is the set of elements of \mathbb{H}_d^n having A_k occurrences of k for each $k = 0, 1, \dots, d$. So

$$|F_1| = \frac{n!}{A_0! A_1! \dots A_d!}. \tag{3.3}$$

For our applications we will be interested in asymptotic behaviour as $n \rightarrow \infty$. To this end, we define W_d as the set of $D = (D_\mu)_{\mu \in \Omega_d}$ such that

$$D_\mu \geq 0 \quad \text{and} \quad \sum_{\mu \in \Omega_d} D_\mu = 1.$$

We regard W_d as a subset of $\mathbb{R}^{m(d)}$, where $m(d) = |\Omega_d|$, and define $Q_i: \mathbb{R}^{m(d)} \rightarrow \mathbb{R}^{d+1}$ in the same way as above.

Now fix $E \in W_d$ and suppose that $\Gamma^{(n)} \in W_d^{(n)}$ with $n^{-1} \Gamma_\mu^{(n)} \rightarrow E_\mu$ for each $\mu \in \Omega_d$. We assume $\Gamma^{(n)}$ (and, hence, E) symmetric. Corresponding to each $\Gamma^{(n)}$, we define $S_0^{(n)}$ and $S^{(n)}$ as above. We then have that $n^{-1} A_k^{(n)} \rightarrow A_k$, where $Q_i E = (A_0, \dots, A_d)$. Then, from (3.1), using (3.2) and (3.3), we deduce

$$\liminf |T^{(n)}|^{1/n} \geq \prod_{k=0}^d A_k^{-A_k} \inf_{D \in \Lambda_E} \prod_{\mu \in \Omega_d} (D_\mu^{D_\mu} / E_\mu^{E_\mu}) \tag{3.4}$$

as $n \rightarrow \infty$, where $\Lambda_E = \{D \in W_d: Q_i D = Q_i E\}$ as before. Note that, as E is assumed to be symmetric, the convex set Λ_E is invariant under permutations of $\{1, 2, 3\}$, and as the logarithm of the product on the right-hand side of (3.4) is convex as a function of D , the infimum is attained at a symmetric D , so it is sufficient to restrict the infimum to symmetric D in Λ_E .

In the type 1 applications d will be 2, 4 or 8. We will find that when $d = 2$, we always have $S = S_0$, and when $d = 4$ we do not have $S = S_0$ but E is the only symmetric element of Λ_E , so the right-hand side of (3.4) is 1. Only in the case $d = 8$ is the symmetric part of Λ_E non-trivial.

3.2. Type 2

In this case we start by fixing $\mu \in \Omega_{2d}$, and define $\Omega_\mu = \{\nu \in \Omega_d: \mu - \nu \in \Omega_d\}$, so that, for example, if $d = 2$ and $\mu = (1, 1, 2)$, then

$$\Omega_\mu = \{(1, 0, 1), (0, 1, 1), (0, 0, 2), (1, 1, 0)\}.$$

Similarly to case 1, we define $W_\mu^{(n)}$ to be the set of $\Gamma = (\Gamma_\nu)_{\nu \in \Omega_\mu}$ with the Γ_ν being non-negative integers with $\sum_{\nu \in \Omega_\mu} \Gamma_\nu = n$. Then given $\Gamma \in W_\mu^{(n)}$, we define \mathbb{S}_0 to be the set of elements of Ω_μ^n containing Γ_ν occurrences of ν for each $\nu \in \Omega_\mu$. And we define $\mathcal{F}_i = P_i(\mathbb{S}_0)$ and $\mathbb{S} = \Omega_\mu^n \cap (\bigcap_{i=1}^3 P_i^{-1}(\mathcal{F}_i))$. Also note that if we write $Q_i \Gamma = (A_{i0}, \dots, A_{id})$, then we have

$$|\mathcal{F}_i| = \frac{n!}{A_{i0}! \cdots A_{id}!}.$$

Then, as in type 1, we can define W_μ and consider a sequence $\Gamma^{(n)} \in W_\mu^{(n)}$ with $\Gamma^{(n)}/n \rightarrow E \in W_\mu$ and then we obtain

$$\left(\frac{|\mathbb{S}^{(n)}|}{|\mathbb{S}_0^{(n)}|}\right)^{1/n} \rightarrow \sup_{D \in \Lambda_{\mu,E}} \prod_{\nu \in \Omega_\mu} \frac{E_\nu^{E_\nu}}{D_\nu^{D_\nu}}, \tag{3.5}$$

where $\Lambda_{\mu,E}$ is the set of $D \in W_\mu$ such that $Q_i D = Q_i E$ for $i = 1, 2, 3$.

We also consider a symmetry property for type 2. We say that $\Gamma \in W_\mu^{(n)}$ or W_μ is symmetric if $\Gamma_{\mu-\nu} = \Gamma_\nu$ for all $\nu \in \Omega_\mu$, and if Γ is unchanged by any permutation that preserves μ . Then, if we assume that $\Gamma^{(n)}$ and E are symmetric, then, as with type 1, in (3.5) we can restrict the supremum to symmetric D .

Type 2 differs from type 1 in the lack of symmetry under permutations of $(1, 2, 3)$. In particular, the sets $\mathcal{F}_1, \mathcal{F}_2$ and \mathcal{F}_3 will in general have different sizes. In order to apply lemma 3.3 effectively, we need to symmetrize. We let π be a cyclic permutation $\pi(\mu) = (\mu_3, \mu_1, \mu_2)$ and then we define $\tilde{\Omega}_\mu = \Omega_\mu \times \Omega_{\pi(\mu)} \times \Omega_{\pi^2(\mu)}$. Now set

$$\begin{aligned} S_0 &= \mathbb{S}_0 \times \pi(\mathbb{S}_0) \times \pi^2(\mathbb{S}_0), \\ F_i &= \mathcal{F}_i \times \pi(\mathcal{F}_i) \times \pi^2(\mathcal{F}_i), \\ S &= \tilde{\Omega}_\mu^n \cap \left(\bigcap_{i=1}^3 P_i^{-1}(F_i)\right) = \mathbb{S} \times \pi(\mathbb{S}) \times \pi^2(\mathbb{S}). \end{aligned}$$

We see that the sets $\mathbb{S}_0, \pi(\mathbb{S}_0)$ and $\pi^2(\mathbb{S}_0)$ all have the same size, and likewise for $\mathbb{S}, \pi(\mathbb{S})$ and $\pi^2(\mathbb{S})$, and so $|S|/|S_0| = (|S|/|\mathbb{S}_0|)^3$. We also have $|F_i| = |\mathcal{F}_1| |\mathcal{F}_2| |\mathcal{F}_3|$

for each i . Then, in the $n \rightarrow \infty$ limit, if we write $Q_i E = (\mathbb{A}_{i0}, \dots, \mathbb{A}_{id})$, we have

$$\liminf |T^{(n)}|^{1/n} \geq \prod_{i=1}^3 \prod_{j=0}^d \mathbb{A}_{ij}^{-\mathbb{A}_{ij}} \inf_{D \in \Lambda_{\mu,E}^*} \prod_{\nu \in \Omega_\mu} (D_\nu^{D_\nu} / E_\nu^{E_\nu}), \tag{3.6}$$

where $\Lambda_{\mu,E}^*$ is the set of symmetric elements of $\Lambda_{\mu,E}$.

We note that, in most of our applications, $\Lambda_{\mu,E}^*$ will be trivial, i.e. equal to $\{E\}$, and then the right-hand side of (3.6) reduces to

$$\prod_{i=1}^3 \prod_{j=0}^d \mathbb{A}_{ij}^{-\mathbb{A}_{ij}}.$$

In the applications we typically have a maximization problem to make the best choice of E , and in the type 2 case the following elementary result will simplify the calculations.

LEMMA 3.4. *Suppose that $k \in \mathbb{N}$ and $\alpha_i > 0$ for $i = 1, \dots, k$. Then the maximum value of $f(x) = \prod_{i=1}^k (\alpha_i/x_i)^{x_i}$ for $x \in \mathbb{R}^k$ subject to $x_i \geq 0$ and $\sum_{i=1}^k x_k = 1$ is $A = \sum_{i=1}^k \alpha_i$, attained for $x_i = \alpha_i/A$.*

We also make the related observation that, more generally, if we have

$$f(x) = \prod_{i=1}^k \alpha_i^{x_i} \prod_{j=1}^l \phi_j(x)^{-\phi_j(x)},$$

where the ϕ_j are affine and real-valued on \mathbb{R}^k . Furthermore, if Z is a convex subset of \mathbb{R}^k on which $x_i \geq 0$ and $\phi_j(x) \geq 0$, then $\log f(x)$ is concave on Z , and so if $y \in Z$ is a critical point of $\log f$, then the maximum of f over Z will be attained at y .

A final note on lemma 3.3 is that a similar combinatorial result is theorem 15.39 of [5]. The proof of this theorem uses a similar probabilistic argument to that given here, but avoids the use of the Salem–Spencer theorem. However, the combinatorial conclusion is somewhat weaker. In [5], the application to the complexity problem makes use of the concept of *degeneration*, which is an elaboration of the idea of border rank, and this enables the weaker combinatorial result to suffice.

4. The Coppersmith–Winograd bounds

Recall the trilinear form χ given by (2.5). The subscripts of the six terms in (2.5) are just the six elements of Ω_2 , in the notation of § 3. Then, given $N \in \mathbb{N}$, the tensor power χ^{3N} can be written as the sum of 6^{3N} terms, each being a tensor product like $\phi_{110} \otimes \phi_{020} \otimes \phi_{011} \cdots$, the sequence of subscripts corresponding to an element of Ω_2^{3N} .

We now follow the discussion of type 1 in § 3, with $d = 2$ and $n = 3N$. We fix a symmetric $\Gamma \in W_2^{3N}$, noting that such a Γ corresponds to a choice of non-negative integers α and β with $\alpha + \beta = N$ and $\Gamma_{200} = \Gamma_{020} = \Gamma_{002} = \alpha$, $\Gamma_{011} = \Gamma_{101} = \Gamma_{110} = \beta$. We define S_0, Q_i, F_i and S as in § 3. We find that

$$Q_1 \Delta = (\Delta_{011} + \Delta_{002} + \Delta_{020}, \Delta_{101} + \Delta_{110}, \Delta_{200}),$$

etc., and it is clear that Δ is determined if we know each $Q_i\Delta$, so $S = S_0$. We also have $A_0 = 2\alpha + \beta$, $A_1 = 2\beta$ and $A_2 = \alpha$.

Now we apply lemma 3.3. Let G_1, G_2, G_3 and T be as given by this lemma, and set to zero all x -variables except those labelled by members of G_1 , and similarly for y and z variables. The resulting trilinear form is a direct sum of $|T|$ copies of \mathcal{M}_{q^β} . Thus, $V_\chi(\rho)^{3N} \geq |T|q^{\beta\rho}$. Now fix $b \in (0, 1)$ and let $N \rightarrow \infty$ with $\beta/N \rightarrow b$ and, using (3.4), we obtain

$$V_\chi(\rho)^3 \geq \frac{27q^{b\rho}}{(1-b)^{1-b}(2b)^{2b}(2-b)^{2-b}}.$$

We choose b to maximize the expression on the right and obtain

$$V_\chi(\rho) \geq \frac{3}{8}\{z^{1/2}(z+32)^{3/2} - z^2 + 80z + 128\}^{1/3}, \tag{4.1}$$

where $z = q^\rho$. Combining this with the bound $\bar{R}(\chi) \leq q + 2$, we obtain an upper bound for ω . The choice $q = 6$ gives the best bound, $\omega \leq \log_6 z$, where z is the unique real root of the cubic equation

$$729z^3 - 64044z^2 + 4458672z - 261404224 = 0.$$

Numerically, we find $z = 72.0435014$ giving $\omega < 2.38719$ as in [7, § 7].

Coppersmith and Winograd improve this bound by considering the tensor square $\chi \otimes \chi$. By taking the tensor product of (2.5) with itself, we can express $\chi \otimes \chi$ as a sum of 36 terms like $\phi_{011} \otimes \phi_{011}$, etc. The improvement is attained by combining some of these terms to make a larger matrix product. We write

$$\phi_{022} = \phi_{011} \otimes \phi_{011} + \phi_{002} \otimes \phi_{020} + \phi_{020} \otimes \phi_{002},$$

and similarly for ϕ_{202} and ϕ_{220} . Likewise, we write

$$\begin{aligned} \phi_{013} &= \phi_{011} \otimes \phi_{002} + \phi_{002} \otimes \phi_{011}, \\ \phi_{004} &= \phi_{002} \otimes \phi_{002}, \\ \phi_{112} &= \phi_{101} \otimes \phi_{011} + \phi_{011} \otimes \phi_{101} + \phi_{002} \otimes \phi_{110} + \phi_{110} \otimes \phi_{002}, \end{aligned}$$

and similarly for ϕ_{103} , etc., obtained by permuting x, y and z . Then we have

$$\chi \otimes \chi = \phi_{004} + \phi_{013} + \phi_{022} + \phi_{112} + \dots, \tag{4.2}$$

where the dots denote the terms obtained from these four by permutation of the subscripts. There are 15 terms in all, one for each of the 15 elements of Ω_4 .

Using a natural notation, we can write

$$\phi_{022}(x, y, z) = x^{[00]} \left(\sum_{i=1}^q \sum_{j=1}^q y_{ij}^{[11]} z_{ij}^{[11]} + y^{[20]} z^{[02]} + y^{[02]} z^{[20]} \right),$$

which is a $(1, 1, q^2 + 2)$ matrix product. Similarly ϕ_{013} is a $(1, 1, 2q)$ matrix product and ϕ_{004} is a $(1, 1, 1)$ matrix product. However, ϕ_{112} is not a matrix product; in order to apply a similar argument to the proof of (4.1) we need a lower bound for the value of ϕ_{112} , which we obtain now. We write V_{112} for $V_{\phi_{112}}$.

LEMMA 4.1. $V_{112}(\rho)^3 \geq 4q^\rho(q^\rho + 2)$.

Proof. We follow the notation and analysis of type 2 in §3, with $d = 2$, $n = 2N$ and $\mu = 112$. We write

$$\phi_{112} = \phi_{011} \otimes \phi_{101} + \phi_{101} \otimes \phi_{011} + \phi_{002} \otimes \phi_{110} + \phi_{110} \otimes \phi_{002}.$$

The four terms in this sum can be labelled by the first of the two subscripts, namely, 011, 101, 002 and 110 in the above order. These are exactly the four elements of Ω_{112} . Then ϕ_{112}^{2N} is the sum of 4^{2N} terms, which can be labelled by the elements of Ω_{112}^{2N} .

Now a symmetric element of W_{112}^{2N} corresponds to a choice of non-negative integers α and β with $\alpha + \beta = N$, and $\Gamma_{002} = \Gamma_{110} = \alpha$, $\Gamma_{101} = \Gamma_{011} = \beta$. We then note that $Q_1(\Gamma) = Q_2(\Gamma) = (N, N, 0)$ and $Q_3(\Gamma) = (\alpha, 2\beta, \alpha)$. Since α and β are determined by Q_3 , we see that $A_{\mu, \Gamma}^*$ is trivial. We fix such a choice of α and β and define \mathbb{S}_0 , \mathcal{F}_i and \mathbb{S} as in type 2.

Now consider $\phi_{112}^{2N} \otimes \phi_{211}^{2N} \otimes \phi_{121}^{2N}$. This can be expressed as a sum of 4^{6N} terms, each labelled by an element of Ω_{112} . We define S_0 , F_i and S as in type 2. Each term labelled by an element of S_0 will be an (m, m, m) matrix product with $m = q^{2\alpha+4\beta}$.

Now apply lemma 3.3, obtaining G_i and T , and set to zero all x -blocks indexed by sequences not in G_1 , and similarly for y - and z -blocks. We obtain a direct sum of $|T|$ matrix products of size (m, m, m) . This implies

$$V_{\phi_{112}}(\rho)^{6N} \geq q^{(2\alpha+4\beta)\rho|T|},$$

and this holds for any choice of positive integers α and β with $\alpha + \beta = N$.

Next, fix a with $0 < a < 1$, and let $N, \alpha \rightarrow \infty$ with $\alpha/N \rightarrow a$. Using (3.6), we obtain

$$V_{\phi_{112}}(\rho)^3 \geq \frac{q^{(2-a)\rho 2^{2+a}}}{a^a(1-a)^{1-a}}.$$

Choosing a to maximize the right-hand side using lemma 3.4 gives $V_{\phi_{112}}(\rho)^3 \geq 4q^\rho(q^\rho + 2)$. \square

Using (4.2), $\chi^{6N} = (\chi \otimes \chi)^{3N}$ can be written as the sum of 15^{3N} trilinear forms, each being a tensor product of $3N$ terms from the set $\{\phi_{004}, \phi_{013}, \phi_{022}, \phi_{112}, \dots\}$, each such form being labelled by an element of Ω_4^{3N} . Again we follow the discussion of type 1, taking $d = 4$ and $n = 3N$. We fix a symmetric $\Gamma \in W_4^{3N}$; such a Γ corresponds to a choice of non-negative integers α, β, γ and δ with $\alpha + 2\beta + \gamma + \delta = N$ such that $\Gamma_{004} = \alpha$, $\Gamma_{013} = \beta$, $\Gamma_{022} = \gamma$ and $\Gamma_{112} = \delta$.

We then find that $Q\Gamma = (2\alpha + 2\beta + \gamma, 2\beta + 2\delta, 2\gamma + \delta, 2\beta, \alpha)$. The linear mapping defined by this expression is clearly injective, so the only symmetric member of A_Γ is Γ itself.

Now fix $\rho \in [2, 3]$ and $\epsilon > 0$, and apply lemma 3.3, obtaining G_1, G_2 and G_3 as before. Again we set all x -variables to zero except those labelled by members of G_1 , and similarly for y and z variables. The resulting trilinear form is a direct sum of $|T|$ copies of $\mathcal{M}_l \otimes (\phi_{112} \otimes \phi_{121} \otimes \phi_{211})^\delta$, where $l = (2q)^{2\beta}(q^2 + 2)^\gamma$. Now, from the bound for V_{112} , we have that $(\phi_{112} \otimes \phi_{121} \otimes \phi_{211})^\delta$ has a restriction that is the direct sum of k copies of \mathcal{M}_n , where

$$kn^\rho \geq C_1 \{4q^\rho(q^\rho + 2)\}^{(1-\epsilon)\delta},$$

where $C_1 > 0$ is independent of N .

Then χ^{6N} has a restriction isomorphic to the direct sum of $k|T|$ copies of \mathcal{M}_{ln} , and so

$$V_\chi(\rho)^{6N} \geq k|T|(ln)^\rho \geq C_1|T|(2q)^{2\beta\rho}(q^2 + 2)^{\gamma\rho}\{4q^\rho(q^\rho + 2)\}^{(1-\epsilon)\delta}.$$

Now fix positive a, b, c and d with $3(a + 2b + c + d) = 1$ and let $N, \alpha, \beta, \gamma, \delta \rightarrow \infty$ with $\alpha/N \rightarrow 3a, \beta/N \rightarrow 3b, \gamma/N \rightarrow 3c$ and $\delta/N \rightarrow 3d$. Then let $\epsilon \rightarrow 0$. Using (3.4), we obtain

$$V_\chi(\rho)^2 \geq \frac{(2q)^{2b\rho}(q^2 + 2)^{c\rho}\{4q^\rho(q^\rho + 2)\}^d}{(2a + 2b + c)^{2a+2b+c}(2b + 2d)^{2b+2d}(2c + d)^{2c+d}(2b)^{2b}a^a}. \tag{4.3}$$

Combined with the bound $V_\chi(\omega) \leq \bar{R}(\chi) \leq q + 2$, this gives a bound for ω , which, in [7, §8], is found numerically to be $\omega < 2.375477$ for the optimal choice of a, b, c and d and $q = 6$.

5. Improved bound

We start with our expression for χ^2 as a sum of 15 terms and take the tensor square, expressing χ^4 as a sum of 225 tensor products. Then, as before, we collect these into groups according to the sums of the x, y and z indices. For example,

$$\begin{aligned} \phi_{125} = \phi_{004} \otimes \phi_{121} + \phi_{013} \otimes \phi_{112} + \phi_{022} \otimes \phi_{103} + \phi_{103} \otimes \phi_{022} \\ + \phi_{112} \otimes \phi_{013} + \phi_{121} \otimes \phi_{004}, \end{aligned}$$

where χ^4 is the sum of 45 such trilinear forms, which can be divided into 10 classes under permutation of the indices, represented by $\phi_{008}, \phi_{017}, \phi_{026}, \phi_{035}, \phi_{044}, \phi_{116}, \phi_{125}, \phi_{134}, \phi_{224}$ and ϕ_{233} .

The first five of these forms, with an x -index of zero, are all of the form $x^{[0000]}$ times a scalar product of y and z vectors. For example,

$$\phi_{026} = \phi_{022} \otimes \phi_{004} + \phi_{013} \otimes \phi_{013} + \phi_{004} \otimes \phi_{022},$$

and the y and z variables occurring in these three tensor products are disjoint, so we have scalar product of size $q^2 + 2 + (2q)^2 + q^2 + 2 = 6q^2 + 4$, i.e. a $(1, 1, n)$ matrix product with $n = 6q^2 + 4$. We get the same for the other four forms, with $n = 1$ for ϕ_{008} , $4q$ for ϕ_{017} , $4q(q^2 + 3)$ for ϕ_{035} and $q^4 + 12q^2 + 6$ for ϕ_{044} .

The other five forms are not matrix products and we need value bounds, as for ϕ_{112} . In order to state these, we introduce some notation. We define quantities E, H and L (depending on q and ρ) as follows: $E = (2q)^\rho, H = (q^2 + 2)^\rho$ and $L = 4q^\rho(q^\rho + 2)$. Then we have the following.

LEMMA 5.1. *If $\rho \in [2, 3]$ and q is an integer with $q > 1$, then*

- (i) $V_{116}(\rho)^3 \geq 4(E^2 + 2L),$
- (ii) $V_{125}(\rho)^3 \geq 4(L + EH)(2H + L)/H,$
- (iii) $V_{134}(\rho)^3 \geq 4(E + L)(2 + 2E + H),$
- (iv) $V_{224}(\rho)^3 \geq (2H + L)^2(2 + 2E + H)/H,$
- (v) $V_{233}(\rho)^3 \geq 4(E + L)^2(2H + L)/L.$

Proof. We first note some inequalities involving E , H and L which we shall need:

$$16 \leq E < H < L < 4H. \quad (5.1)$$

The first two are immediate since $q \geq 2$ and $\rho \geq 2$, and we also have $(q^2+2)/q^2 \leq \frac{3}{2}$, so, using $\rho \leq 3$, we deduce $H \leq (\frac{3}{2})^3 q^{2\rho} < L$. Furthermore,

$$(q^2 + 2)^\rho \geq q^{2\rho} + 2\rho q^{2(\rho-1)} > q^{2\rho} + 2q^\rho,$$

so $L < 4H$.

Now we prove the individual bounds.

(i) We write

$$\phi_{116} = \phi_{013} \otimes \phi_{103} + \phi_{103} \otimes \phi_{013} + \phi_{004} \otimes \phi_{112} + \phi_{112} \otimes \phi_{004}.$$

This is the same as for ϕ_{112} , except that the third indices have all increased by 2. Again we follow type 2, with $d = 2$, $n = 2N$ and $\mu = 116$. As for ϕ_{112} , a symmetric element of W_{116}^{2N} corresponds to a choice of $\alpha, \beta \geq 0$ with $\alpha + \beta = N$ and $\Gamma_{004} = \alpha$, etc. As before, we can express $\phi_{116}^{2N} \otimes \phi_{611}^{2N} \otimes \phi_{161}^{2N}$ as a sum of 4^{6N} terms, each labelled by an element of $\tilde{\Omega}_{116}$. The combinatorics of these sequences is exactly the same as for 112, and we can use the same G_i and T as in lemma 4.1.

Then each member of T indexes a trilinear form isomorphic to $\mathcal{M}_l \otimes (\phi_{112} \otimes \phi_{121} \otimes \phi_{211})^\alpha$, where $l = (2q)^{2\beta}$. We then proceed as in the last part of the proof of (4.3), and obtain $V_{116}(\rho)^{6N} \geq k|T|(lm)^\rho$, where $km^\rho \geq CL^{(1-\epsilon)\alpha}$. Using the bound for $|T|$ from the proof of lemma 4.1, and letting $N, \alpha \rightarrow \infty$ with $\alpha/N \rightarrow a$, we then get

$$V_{116}(\rho)^3 \geq \frac{2(2L)^a E^{2(1-a)}}{a^a(1-a)^{1-a}}$$

for $0 < a < 1$ and, by optimizing a using lemma 3.4, we deduce $V_{116}(\rho)^3 \geq 4(2L + E^2)$.

(ii) Again, we write ϕ_{125} as a sum of six tensor products, and $\phi_{125}^{2N} \otimes \phi_{512}^{2N} \otimes \phi_{251}^{2N}$ as a sum of 6^{6N} terms, each labelled by an element of $\tilde{\Omega}_{125}^{2N}$. A symmetric element Γ of $W_{125}^{(2N)}$ corresponds to a choice of $\alpha, \beta, \gamma \geq 0$ with sum N such that $\Gamma_{004} = \Gamma_{121} = \alpha$, $\Gamma_{013} = \Gamma_{112} = \beta$ and $\Gamma_{022} = \Gamma_{103} = \gamma$. Then

$$\begin{aligned} Q_1\Gamma &= (N, N, 0, 0, 0), \\ Q_2\Gamma &= (\alpha + \gamma, 2\beta, \alpha + \gamma, 0, 0), \\ Q_3\Gamma &= (0, \alpha, \beta + \gamma, \beta + \gamma, \alpha). \end{aligned}$$

Again, Q_2 and Q_3 together determine α, β and γ so $\Lambda_{125, \Gamma}^*$ is trivial.

As before, we obtain G_i and T , and after ‘setting to zero’ we have a direct sum of $|T|$ copies of

$$\mathcal{M}_l \otimes (\phi_{112} \otimes \phi_{121} \otimes \phi_{211})^{\alpha+\beta},$$

where $l = (2q)^{\beta+\gamma}(q^2+2)^\gamma$. Taking limits with $\alpha/N \rightarrow a$ and $\beta/N \rightarrow b$, and using (3.6), we find that, for $a, b \geq 0$ with $a + b \leq 1$, we have

$$V_{125}^3 \geq 4H^{-1} \left\{ \frac{L^a (EH)^{1-a}}{a^a(1-a)^{1-a}} \right\} \left\{ \frac{L^b (2H)^{1-b}}{b^b(1-b)^{1-b}} \right\}.$$

Using lemma 3.4, the two terms in curly brackets can be maximized separately, when $a = L/(L + EH)$ and $b = L/(L + 2H)$. We need to check that these values satisfy $a + b \leq 1$; this is equivalent to $L^2 \leq 2EH^2$, which follows from the first and last inequalities in (5.1). Using these a and b we obtain the required bound for $V_{125}(\rho)$.

(iii) We proceed in the same way as before. A symmetric $\Gamma \in W_{134}^{(2N)}$ corresponds to $\alpha, \beta, \gamma, \delta \geq 0$ with sum N such that $\Gamma_{004} = \alpha, \Gamma_{013} = \beta, \Gamma_{022} = \gamma$ and $\Gamma_{031} = \delta$. We then have

$$\begin{aligned} Q_1\Gamma &= (N, N, 0, 0, 0), \\ Q_2\Gamma &= (\alpha + \delta, \beta + \gamma, \beta + \gamma, \alpha + \delta, 0), \\ Q_3\Gamma &= (\alpha, \beta + \delta, 2\gamma, \beta + \delta, \alpha), \end{aligned}$$

and $\Lambda_{134, \Gamma}^*$ is again trivial.

This time we get a direct sum of $|T|$ copies of

$$\mathcal{M}_l \otimes (\phi_{112} \otimes \phi_{121} \otimes \phi_{211})^{\beta+\gamma},$$

where $l = (2q)^{\alpha+\beta+2\delta}(q^2+2)^\gamma$. We find, by taking the limit with $\alpha/N \rightarrow a, \beta/N \rightarrow b$ and $\gamma/N \rightarrow c$ that if $a, b, c \geq 0$ with $a + b + c \leq 1$, then

$$V_{134}(\rho)^3 \geq \frac{2^{3-c}E^{2-a-b-2c}H^cL^{b+c}}{(b+c)^{b+c}(1-b-c)^{1-b-c}a^ac^c(1-a-c)^{1-a-c}}.$$

If we put $\sigma = b + c$, we can write this as

$$V_{134}(\rho)^3 \geq 8 \left\{ \frac{E^{1-\sigma}L^\sigma}{\sigma^\sigma(1-\sigma)^{1-\sigma}} \right\} \left\{ \frac{E^{1-a-c}(H/2)^c}{a^ac^c(1-a-c)^{1-a-c}} \right\}.$$

Using lemma 3.4 as before, we can then maximize the expressions in curly brackets by taking $\sigma = L/(E + L)$, $a = 2/(2 + 2E + H)$ and $c = H/(2 + 2E + H)$ and obtain the required bound, but we need to check that the conditions on a, b and c are satisfied, and for that we require $\sigma \geq c$ and $\sigma + a \leq 1$. These are respectively equivalent to $EH \leq 2L(1 + E)$ and $2L \leq EH + 2E^2$, the first of which follows from $H < L$ and the second from $E \geq 16$ and $L < 4H$ in (5.1).

(iv) This time a symmetric $\Gamma \in W_{224}^{(2N)}$ corresponds to $\alpha, \beta, \gamma, \delta \geq 0$ such that $\alpha + 2\beta + \gamma + \delta = N$ and $\Gamma_{004} = \alpha, \Gamma_{013} = \beta, \Gamma_{022} = \gamma$ and $\Gamma_{112} = \delta$. We then have

$$Q_1\Gamma = Q_2\Gamma = (\alpha + \beta + \gamma, 2\beta + 2\delta, \alpha + \beta + \gamma, 0, 0)$$

and $Q_3\Gamma = (\alpha, 2\beta, 2\gamma + 2\delta, 2\beta, \alpha)$; again $\Lambda_{224, \Gamma}^*$ is trivial.

Then we get a direct sum of $|T|$ copies of $\mathcal{M}_l \otimes (\phi_{112} \otimes \phi_{121} \otimes \phi_{211})^{2\beta+2\delta}$, where $l = (2q)^{2\beta}(q^2 + 2)^{\alpha+2\gamma}$. We find, by taking the limit with $\alpha/N \rightarrow a, \beta/N \rightarrow b$ and $\gamma/N \rightarrow c$, that if $a, b, c \geq 0$ with $a + 2b + c \leq 1$ then

$$V_{224}(\rho)^3 \geq \frac{2^{3a+4b+2c}E^{2b}H^{a+2c}L^{2(1-a-b-c)}}{\{(a+b+c)^{a+b+c}(1-a-b-c)^{1-a-b-c}\}^2 a^a(2b)^{2b}(1-a-2b)^{1-a-2b}}.$$

If we put $\sigma = a + b + c$ and $b' = 2b$ we can write this as

$$V_{224}(\rho)^3 \geq \left\{ \frac{L^{1-\sigma}(2H)^\sigma}{\sigma^\sigma(1-\sigma)^{1-\sigma}} \right\}^2 \left\{ \frac{(2/H)^a(2E/H)^{b'}}{a^ab'^{b'}(1-a-b')^{1-a-b'}} \right\}.$$

By lemma 3.4 again we can maximize the expressions in curly brackets by taking $\sigma = 2H/(2H + L)$, $a = 2/(2 + 2E + H)$ and $b = E/(2 + 2E + H)$ and obtain the required bound, but we need to check that the conditions on a , b and c are satisfied. For this, we require $a + b \leq \sigma \leq 1 - b$. This is equivalent to $(2 + E)L \leq 2H(E + H)$ together with $2EH \leq L(2 + E + H)$, which both follow from (5.1).

(v) This time a symmetric $\Gamma \in W_{233}^{(2N)}$ corresponds to $\alpha, \beta, \gamma, \delta \geq 0$ such that $2\alpha + \beta + \gamma + \delta = N$ and $\Gamma_{013} = \alpha$, $\Gamma_{022} = \beta$, $\Gamma_{103} = \gamma$ and $\Gamma_{112} = \delta$. We then have

$$\begin{aligned} Q_1\Gamma &= (2\alpha + \beta, 2\gamma + 2\delta, 2\alpha + \beta, 0, 0), \\ Q_2\Gamma &= Q_3\Gamma = (\alpha + \gamma, \alpha + \beta + \delta, \alpha + \beta + \delta, \alpha + \gamma, 0). \end{aligned}$$

In this case it is not true that $\Lambda_{233,\Gamma}^*$ is trivial since $Q_i\Gamma$ are determined if we know $2\alpha + \beta$ and $\alpha + \gamma$. This has to be taken into account in the estimation of $|T|$.

Suppose $\sigma, \mu \geq 0$ with $\sigma + 2\mu \leq 2$. These conditions ensure that we can find $a, b, c \geq 0$ such that $2a + b + c \leq 1$ and $\sigma = 2a + b$ and $\mu = a + c$. We suppose that a, b and c are chosen to minimize $a^{2a}b^bc^c(1 - 2a - b - c)^{1-2a-b-c}$ subject to these constraints. This will ensure that the infimum in (3.6) is 1. Then we obtain a direct sum of $|T|$ copies of

$$\mathcal{M}_l \otimes (\phi_{112} \otimes \phi_{121} \otimes \phi_{211})^{\beta+2\delta},$$

where $l = (2q)^{2\alpha+2\gamma}(q^2 + 2)^{2\alpha+\beta}$. Then we find, by taking the limit with $\alpha/N \rightarrow a$, $\beta/N \rightarrow b$ and $\gamma/N \rightarrow c$, that

$$V_{233}(\rho)^3 \geq \frac{2^{2+2a+b}E^{2a+2c}H^{2a+b}L^{2-4a-b-2c}}{(2a + b)^{2a+b}(1 - 2a - b)^{1-2a-b}\{(a + c)^{a+c}(1 - a - c)^{1-a-c}\}^2}.$$

Conveniently, the right-hand side can be expressed in terms of σ and μ , and we have

$$V_{233}(\rho)^3 \geq 4L^2 \left\{ \frac{(2H/L)^\sigma}{\sigma^\sigma(1 - \sigma)^{1-\sigma}} \right\} \left\{ \frac{(E/L)^\mu}{\mu^\mu(1 - \mu)^{1-\mu}} \right\}^2.$$

We can then maximize the expressions in curly brackets by taking $\sigma = 2H/(2H + L)$ and $\mu = E/(E + L)$ and obtain the required bound, but we need to check that this choice of σ and μ satisfies $\sigma + 2\mu \leq 2$. This is equivalent to $EH \leq L(H + L)$, which follows from $E < L$ in (5.1). \square

Table 1 lists the 10 classes of trilinear form, indexed by $i = 1, \dots, 10$, and the cubes, denoted by v_i , of the lower bounds we have obtained for their values.

The table also shows the number n_i such that there are $3n_i$ forms in the i th class. For $i = 1, \dots, 5$, we also write u_i for the actual matrix size, so that $u_2 = 4q$, etc., and then $v_i = u_i^3$.

Using our expression for χ^4 as a sum of 45 forms, we can write $\chi^{12N} = (\chi^4)^{3N}$ as a sum of 45^{3N} trilinear forms, each being a tensor product of $3N$ terms from the set $\{\phi_{008}, \phi_{017}, \dots\}$ of 45 forms, each such tensor product being labelled by an element of $\Omega_{8,3N}$. Again, we follow the type 1 analysis, now with $d = 8$ and $n = 3N$. Let Z denote the set of $a = (a_1, \dots, a_{10}) \in \mathbb{R}^{10}$ such that each $a_i \geq 0$ and $\sum_{i=1}^{10} n_i a_i = 1$ and let Z_N be the set of $\alpha \in \mathbb{N}^{10}$ such that $\alpha_i \geq 0$ and $\sum_{i=1}^{10} n_i \alpha_i = N$. There is then a one-to-one correspondence between symmetric elements Γ of W_8^{3N} and

Table 1. 10 classes of trilinear form

i	Representative	n_i	v_i
1	ϕ_{008}	1	1
2	ϕ_{017}	2	$(4q)^\rho$
3	ϕ_{026}	2	$(6q^2 + 4)^\rho$
4	ϕ_{035}	2	$\{4q(q^2 + 3)\}^\rho$
5	ϕ_{044}	1	$(q^4 + 12q^2 + 6)^\rho$
6	ϕ_{116}	1	$4(E^2 + 2L)$
7	ϕ_{125}	2	$4(L + EH)(2H + L)/H$
8	ϕ_{134}	2	$4(E + L)(2 + 2E + H)$
9	ϕ_{224}	1	$(2H + L)^2(2 + 2E + H)/H$
10	ϕ_{233}	1	$4(E + L)^2(2H + L)/L$

$\alpha \in Z_N$ given by $\alpha_1 = \Gamma_{008}$, $\alpha_2 = \Gamma_{017}$, etc. Furthermore (for symmetric Γ), we have $Q_1\Gamma = (A_0, \dots, A_8)$, where

$$\left. \begin{aligned} A_0 &= 2\alpha_1 + 2\alpha_2 + 2\alpha_3 + 2\alpha_4 + \alpha_5, \\ A_1 &= 2\alpha_2 + 2\alpha_6 + 2\alpha_7 + 2\alpha_8, \\ A_2 &= 2\alpha_3 + 2\alpha_7 + 2\alpha_9 + \alpha_{10}, \\ A_3 &= 2\alpha_4 + 2\alpha_8 + 2\alpha_{10}, \\ A_4 &= 2\alpha_5 + 2\alpha_8 + \alpha_9, \\ A_5 &= 2\alpha_4 + 2\alpha_7, \\ A_6 &= 2\alpha_3 + \alpha_6, \\ A_7 &= 2\alpha_2, \\ A_8 &= \alpha_1. \end{aligned} \right\} \quad (5.2)$$

We can write (5.2) in matrix form as $A = Q\alpha$, where Q is the 9×10 matrix of coefficients of (5.2). A difference from the earlier case is that the linear mapping defined by Q has non-trivial kernel, as is clear from the dimensions. In fact, if we define $Y = \{x \in \mathbb{R}^{10} : Qx = 0\}$, then Y is a two-dimensional subspace spanned by $\sigma = (0, 0, 1, 0, -2, -2, 0, 2, 0, -2)$ and $\tau = (0, 0, 0, 1, -2, 0, -1, 1, 2, -2)$.

As before, we fix $\rho \in [2, 3]$ and $\epsilon > 0$, choose $\alpha \in Z_N$ (and corresponding symmetric Γ), and apply lemma 3.3, obtaining G_i and T . Now from the definition of V_{116} , there is a constant C independent of N such that $(\phi_{116} \otimes \phi_{161} \otimes \phi_{611})^{\alpha_6}$ has a restriction which is a direct sum of k_6 copies of \mathcal{M}_{m_6} , where $k_6 m_6^\rho \geq C v_6^{(1-\epsilon)\alpha_6}$, and similarly with 6 replaced by 7, 8, 9, and 10. We conclude that χ^{12N} has a restriction isomorphic to the direct sum of $|T| \prod_{j=6}^{10} k_j$ copies of \mathcal{M}_l , where

$$l = \prod_{j=1}^5 u_j^{n_j \alpha_j} \prod_{j=6}^{10} m_j^{n_j \alpha_j}.$$

Then

$$V_\chi(\rho)^{12N} \geq |T| l^\rho \prod_{j=6}^{10} k_j \geq C^5 |T| \prod_{i=1}^{10} v_i^{n_i \alpha_i (1-\epsilon)}.$$

Table 2. Values of a_i and b_i

i	a_i	b_i
1	0.0000001098	0.0000001098
2	0.0000189244	0.0000189244
3	0.0007891925	0.0007968493
4	0.0105046192	0.0106613542
5	0.0371434106	0.0368146264
6	0.0010677681	0.0010524545
7	0.0219568746	0.0218001397
8	0.1409798025	0.1411518518
9	0.2177976006	0.2181110728
10	0.3954922843	0.3951634971

We now proceed to the limit as $N \rightarrow \infty$. The application of (3.4) now requires a minimization, which is handled by the next lemma. We define

$$\mathcal{N} = \{a \in Z : a_1 a_5 a_{10} = a_3^2 a_8 \text{ and } a_2 a_5 a_{10} = a_3 a_4 a_9\}.$$

Then we have the following.

LEMMA 5.2. *Suppose $b \in \mathcal{N}$ and $a \in Z$ with $a - b \in Y$. Then*

$$\prod_{i=1}^{10} b_i^{n_i b_i} \leq \prod_{i=1}^{10} a_i^{n_i a_i}.$$

Proof. Recall the basis σ, τ for Y . Define $f: \mathbb{R}^2 \rightarrow \mathbb{R}^{10}$ by $f(s, t) = b + s\sigma + t\tau$ and define $h: Z \rightarrow \mathbb{R}$ by

$$h(a) = \sum_{i=1}^{10} n_i a_i \log a_i.$$

Then $H(t) = h(f(t))$ is a convex function on the convex set $f^{-1}(W)$, and a simple calculation, using $b \in \mathcal{N}$, shows that the gradient of h vanishes at $(0, 0)$. Hence, h attains its infimum over $f^{-1}(Z)$ at $(0, 0)$ and the result follows. \square

Now suppose that a and b are as in lemma 5.2 and suppose $N^{-1}\alpha \rightarrow a$. Then, applying (3.4) and using lemma 5.2, we conclude the following.

THEOREM 5.3. *Let $b \in \mathcal{N}$ and $a \in Z$ with $a - b \in Y$. Then*

$$V_\chi(\omega)^4 \geq \prod_{i=1}^{10} \{v_i^{a_i/3} a_i^{a_i} b_i^{-b_i}\}^{n_i} \prod_{j=0}^8 \mathbb{A}_j^{-\mathbb{A}_j}, \quad (5.3)$$

where $\mathbb{A} = \frac{1}{3}Qa$.

The bound obtained by applying this theorem with $q = 6$ and the values of a and b in table 2 is $\omega < 2.373689703$.

These values of a and b were obtained by maximizing the right-hand side of (5.3) subject to $b \in \mathcal{N}$ and $a \in Z$ with $a - b \in Y$. This is, in effect, an optimization problem in nine variables, as \mathcal{N} is a seven-dimensional manifold that can be

parametrized by a suitable choice of seven of the variables b_i and then, for a given b , the set of allowed a is a convex subset of two-dimensional affine subspace.

Note added 25 February 2012

Since this paper was submitted, we have learned of independent work by Virginia Williams, who has achieved an improved bound $\omega < 2.3727$ using essentially the same method as in this paper, but using the eighth tensor power of χ rather than the fourth. Her calculations also indicate that, using the fourth power, taking $q = 5$ gives a better bound than we obtained above with $q = 6$.

Acknowledgements

A.J.S. was supported by the EPSRC.

References

- 1 F. A. Behrend. On sets of integers which contain no three terms in arithmetical progression. *Proc. Natl Acad. Sci. USA* **32** (1946), 331–332.
- 2 D. Bini. Relations between exact and approximate bilinear algorithms. *Calcolo* **17** (1980), 87–97.
- 3 D. Bini, M. Capovani, G. Lotti and F. Romani. $O(n^{2.7799})$ complexity for matrix multiplication. *Info. Process. Lett.* **8** (1979), 234–235.
- 4 M. Bläser. A $\frac{5}{2}n^2$ lower bound for the rank of $n \times n$ matrix multiplication over arbitrary fields. *Proceedings of the 40th Annual Symposium on Foundations of Computer Science*, pp. 45–50 (New York: IEEE, 1999).
- 5 P. Bürgisser, M. Clausen and M. A. Shokrollahi. *Algebraic complexity theory*, Comprehensive Studies in Mathematics, volume 315 (Springer, 1997).
- 6 D. Coppersmith and S. Winograd. On the asymptotic complexity of matrix multiplication. *SIAM J. Comput.* **11** (1982), 472–492.
- 7 D. Coppersmith and S. Winograd. Matrix multiplication via arithmetic progressions. *J. Symb. Computat.* **9** (1990), 251–280.
- 8 V. Pan. Strassen’s algorithm is not optimal: trilinear technique of aggregating, uniting and canceling for constructing fast algorithms for matrix operations. *Proceedings of the 19th Annual Symposium on Foundations of Computer Science*, pp. 166–176 (New York: IEEE, 1978).
- 9 V. Pan. New fast algorithms for matrix operations. *SIAM J. Comput.* **9** (1980), 321–342.
- 10 V. Pan. New combinations of methods for the acceleration of matrix multiplication. *Comput. Math. Appl.* **7** (1981), 73–125.
- 11 V. Pan. How can we speed up matrix multiplication? *SIAM Rev.* **26** (1984), 393–415.
- 12 F. Romani. Some properties of disjoint sums of tensors related to matrix multiplication. *SIAM J. Comput.* **11** (1982), 263–267.
- 13 R. Salem and D. C. Spencer. On sets of integers which contain no three terms in arithmetical progression. *Proc. Natl Acad. Sci. USA* **28** (1942), 561–563.
- 14 A. Schönhage. Partial and total matrix multiplication. *SIAM J. Comput.* **10** (1981), 434–455.
- 15 A. Stothers. On the complexity of matrix multiplication. Doctoral Thesis, University of Edinburgh (2010).
- 16 V. Strassen. Gaussian elimination is not optimal. *Numer. Math.* **13** (1969), 354–356.
- 17 V. Strassen. Relative bilinear complexity and matrix multiplication. *J. Reine Angew. Math.* **375** (1987), 406–443.

(Issued 5 April 2013)

