

# SZEMERÉDI'S PROOF OF SZEMERÉDI'S THEOREM

TERENCE TAO

ABSTRACT. We present Szemerédi's original proof of Szemerédi's theorem.

## 1. INTRODUCTION

The purpose of this note is to present Szemerédi's original proof of *Szemerédi's theorem*:

**Theorem 1.1** (Szemerédi's theorem). [3] *Let  $A$  be a subset of  $\mathbf{Z}$  with positive upper density*

$$\limsup_{N \rightarrow \infty} \frac{|A \cap [-N, N]|}{|[-N, N]|} = \delta > 0$$

*and fix  $k \geq 3$ . Then  $A$  contains at least one proper arithmetic progression  $a, a + r, \dots, a + (k - 1)r$  of length  $k$  and  $r > 0$ .*

We have restricted to the case  $k \geq 3$  since the cases  $k < 3$  are trivial. Once one obtains one progression, it is an easy matter to in fact obtain infinitely many progressions (e.g. by deleting the progression that we locate from  $A$  - which does not affect the upper density - and then starting over).

It is well known that Szemerédi's proof relies in such ingredients as Szemerédi's regularity lemma, which essentially allows one to approximate an arbitrary large graph by an object of bounded complexity, as well as van der Waerden's theorem, which allows one to locate a monochromatic arithmetic progression of length  $k$  in any coloring of a sufficiently large arithmetic progression. The proof also relies on another fact, which is that most pairs of integers can be connected by a progression of length  $k$  (specifically, given  $1 \leq j < j' \leq r$  and  $n, m \in \mathbf{Z}$ , there often exists a progression  $a, a + r, \dots, a + (k - 1)r$  with  $a + jr = n$  and  $a + j'r = m$ ). This particular property (needed for the regularity lemma to be useful) is not true for generalizations of Szemerédi's theorem, such as multidimensional, Hales-Jewett, or polynomial analogues, and so it is not clear at this point whether Szemerédi's argument could extend to cover these cases (unless perhaps one replaces the regularity lemma with a hypergraph analogue).

Let us close this introduction with a simple lemma, which basically asserts that Szemerédi's theorem is easy if the density is sufficiently large.

**Lemma 1.2.** *Let  $k \geq 3$ , and let  $P = \{b, b + s, \dots, b + (l - 1)s\}$  be a progression of some length  $l \geq 10k$ . Let  $A$  be a subset of  $P$  such that  $|A| \geq (1 - \frac{1}{10k})|P|$ . Then  $A$  contains at least one proper arithmetic progression  $a, a + r, \dots, a + (k - 1)r$  of length  $k$  with  $r > 0$ .*

**Proof** This is an easy counting argument. The number of progressions of the form  $a, a + r, \dots, a + (k - 1)r$  in  $P$  with  $r > 0$  is basically  $\frac{1}{2k}|P|^2$ , whereas for each  $0 \leq i < k$  and  $b \in P$  the number of such progressions with  $a + ir = b$  is at most  $\frac{1}{k}|P|$ . Thus for each  $i$ , the number of progression with  $a + ir \notin A$  is at most  $\frac{1}{10k^2}|P|^2$ , so summing over all  $i$  we see that there must be at least one progression with  $a + ir \in A$  for all  $0 \leq i < k$ . ■

## 2. PROGRESSIONS

We recall the concept of a generalized arithmetic progression, or *progression* for short.

**Definition 2.1.** If  $N^{(d)} = (N_1, \dots, N_d)$  is a  $d$ -tuple of positive integers for some  $d > 0$ , we form the discrete box

$$[0, N^{(d)}] := [0, N_1] \times \dots \times [0, N_d] = \{(n_1, \dots, n_d) \in \mathbf{Z}^d : 0 \leq n_i < N_i \text{ for all } 1 \leq i \leq d\}.$$

A *progression*  $P$  of rank  $d$  and dimensions  $N^{(d)} = (N_1, \dots, N_d)$  is any subset of  $\mathbf{Z}$  of the form

$$P = a + [0, N^{(d)}] \cdot v = \{a + n \cdot v : n \in [0, N^{(d)}]\} = \{a + \sum_{i=1}^d n_i v_i : 0 \leq n_i < N_i \text{ for all } 1 \leq i \leq d\},$$

where  $v = (v_1, \dots, v_d) \in \mathbf{Z}^d$  is a vector. If the map  $n \mapsto n \cdot v$  is one-to-one on  $[0, N^{(d)})$  we say that  $P$  is *proper* (this is equivalent to the statement that  $|P| = N_1 \dots N_d$ ). For  $0 \leq n_d < N_d$ , we define the  $n_d^{\text{th}}$  component  $P_{n_d}$  of  $P$  to be

$$P_{n_d} = \{a + \sum_{i=1}^d n_i v_i : 0 \leq n_i < N_i \text{ for all } 1 \leq i \leq d - 1\},$$

thus each  $P_{n_d}$  is a progression of rank  $d - 1$  and dimensions  $(N_1, \dots, N_{d-1})$ . We can iterate this definition, defining  $P_{n_d, n_{d-1}}$ , etc.

For Szemerédi's argument we need a couple more notions associated to that of a progression. We first recursively define what it means for a progression to be increasing.

**Definition 2.2.** A progression of rank 0 is always increasing. If  $d > 0$ , a progression  $P$  of rank  $d$  is said to be *increasing* if each component  $P_{n_d}$  is increasing, and if whenever  $0 \leq n_d < n'_d < N_d$ , every element of  $P_{n'_d}$  is larger than every element of  $P_{n_d}$ .

Note that increasing progressions are automatically proper. In our analysis we shall deal almost exclusively with increasing progressions.

Now we define the notion of a *color* of an increasing progression (with respect to a fixed set  $A$ ).

**Definition 2.3.** Let  $2^{[0, N^{(d)})}$  be the power set of  $[0, N^{(d)})$ , i.e. the set of all subsets of  $[0, N^{(d)})$ . The *color*  $c(P) \in 2^{[0, N^{(d)})}$  of an increasing progression  $P = a + [0, N^{(d)}) \cdot v$  is defined to be the set

$$c(P) := \{n \in [0, N^{(d)}) : a + n \cdot v \in A\}.$$

We say that  $P$  is *black* if  $c(P)$  is the empty set.

Note that the color  $c(P)$  of a progression determines the colors  $c(P_0), \dots, c(P_{N_d-1})$  of its components, and conversely. Generally speaking, we wish to avoid the black color whenever we can, and work instead with “saturated” colors which are as “bright as possible” (roughly speaking, they contain as many points as possible; the precise definition is more recursive than this).

To prove Szemerédi’s theorem (Theorem 1.1), our strategy shall be to first work in a progression of large rank and even larger dimensions which contains lots of “saturated” colors, and then trade in this rank and saturation for an increasingly “perfect” progression, ultimately ending up with a progression of length  $k$  which is contained entirely in  $A$ . We begin by describing a certain recursive procedure that will set up for us such a large rank progression.

### 3. ASYMPTOTIC UPPER COLOR DENSITY

Let us fix some dimensions  $N^{(d)} = (N_1, \dots, N_d)$ . The increasing progressions of dimensions  $[0, N^{(d)})$  can take any color within  $2^{[0, N^{(d)})}$ ; let us suppose that we have some subset  $\Sigma_{N^{(d)}} \subset 2^{[0, N^{(d)})}$  of these colors, which we refer to as the *saturated* colors (roughly speaking, they are the “brightest” colors which can occur in progressions of dimensions  $[0, N^{(d)})$ , and will be quite far away from being black). For this discussion the exact definition of saturation is not relevant. Given any increasing progression  $P$  of dimensions  $(N_1, \dots, N_{d+1})$  for some  $N_{d+1}$ , define the *saturation*  $0 \leq \sigma_{N^{(d)}}(P) \leq 1$  of  $P$  to be the quantity

$$\sigma_{N^{(d)}}(P) := \frac{1}{N_{d+1}} |\{0 \leq j < N_{d+1} : c(P_j) \in \Sigma_{N^{(d)}}\}|,$$

i.e.  $\sigma_{N^{(d)}}(P)$  is the proportion of components of  $P$  whose color is saturated. Note that this quantity depends only on the color  $c(P)$  of  $P$  and not on  $P$  itself. We then define the upper saturation numbers

$$\overline{\sigma}_{N^{(d)}}(N_{d+1}) := \sup_P \sigma_{N^{(d)}}(P),$$

where  $P$  ranges over all increasing progressions of dimensions  $(N_1, \dots, N_{d+1})$ , and the *asymptotic upper saturation*

$$\overline{\sigma}_{N^{(d)}}(\infty) := \limsup_{N_{d+1}} \overline{\sigma}_{N^{(d)}}(N_{d+1}).$$

This limit superior is in fact a limit:

**Lemma 3.1.** *If we define the quantity*

$$\mu_d(M) := |\overline{\sigma_{N^{(d)}}}(M) - \overline{\sigma_{N^{(d)}}}(\infty)|,$$

then

$$\lim_{M \rightarrow \infty} \mu_d(M) = 0. \tag{1}$$

**Proof** Let  $\varepsilon > 0$  be arbitrary. It will suffice to show that

$$\overline{\sigma_{N^{(d)}}}(M) \geq \overline{\sigma_{N^{(d)}}}(\infty) - O(\varepsilon)$$

for  $M$  sufficiently large depending on  $\varepsilon$ . But we can find  $M' \gg M$  arbitrarily large such that

$$\overline{\sigma_{N^{(d)}}}(M') = \overline{\sigma_{N^{(d)}}}(\infty) - O(\varepsilon),$$

and thus we can find an increasing progression  $P$  of dimensions  $(N_1, \dots, N_d, M')$  such that

$$\sigma_{N^{(d)}}(P) = \overline{\sigma_{N^{(d)}}}(\infty) - O(\varepsilon).$$

By truncating  $P$  by a little bit we may assume  $M'$  is a multiple of  $M$  (the error term in doing so is negligible if  $M'$  is large enough depending on  $M$ ). But then by subdividing  $P$  into  $M'/M$  blocks of dimension  $(N_1, \dots, N_d, M)$  and using the pigeonhole principle, we can find a progression  $P'$  of dimension  $(N_1, \dots, N_d, M)$  such that

$$\sigma_{N^{(d)}}(P') = \overline{\sigma_{N^{(d)}}}(\infty) - O(\varepsilon).$$

and the claim follows. ■

Let us dispose of an easy case when  $\mu_d$  actually does reach zero:

**Lemma 3.2.** *If  $\overline{\sigma_{N^{(d)}}}(\infty) \geq 1/2$  and  $\mu_d(M) < c(k, N_1, \dots, N_d, M)$  for some sufficiently small  $c(k, N_1, \dots, N_d, M) > 0$  and black is not a saturated color, then we can find a progression of length  $k$  in  $A$ .*

**Proof** Suppose that  $\mu_d(M) < c = c(k, N_1, \dots, N_d, M)$ . By (1), we can find an arbitrarily large  $M'$  such that  $\mu_d(MM') < \mu_d(M)$ . In particular we can find an increasing progression  $P$  of rank  $d+1$  and dimensions  $(N_1, \dots, N_d, MM')$  with

$$\sigma_{N^{(d)}}(P) \geq \overline{\sigma_{N^{(d)}}}(\infty) - \mu_d(M).$$

We now view  $P$  instead as a progression  $Q$  of rank  $d+2$  and dimensions  $(N_1, \dots, N_d, M, M')$ , and observe that

$$\frac{1}{M'} \sum_{n=0}^{M'-1} \sigma_{N^{(d)}}(Q_n) = \sigma_{N^{(d)}}(P) \geq \overline{\sigma_{N^{(d)}}}(\infty) - \mu_d(M)$$

and thus

$$\frac{1}{M'} \sum_{n=0}^{M'-1} (\sigma_{N^{(d)}}(Q_n) - \overline{\sigma_{N^{(d)}}}(\infty)) > -\mu_d(M)$$

On the other hand, we have

$$\sigma_{N^{(d)}}(Q_n) \leq \overline{\sigma_{N^{(d)}}}(M) \leq \overline{\sigma_{N^{(d)}}}(\infty) + \mu_d(M)$$

and hence

$$\frac{1}{M'} \sum_{n=0}^{M'-1} (\sigma_{N^{(d)}}(Q_n) - \overline{\sigma_{N^{(d)}}}(\infty))_+ \leq \mu_d(M)$$

and hence

$$\frac{1}{M'} \sum_{n=0}^{M'-1} |\sigma_{N^{(d)}}(Q_n) - \overline{\sigma_{N^{(d)}}}(\infty)| \leq 3\mu_d(M).$$

By Chebyshev, we see in particular that

$$|\{0 \leq n \leq M' - 1 : |\sigma_{N^{(d)}}(Q_n) - \overline{\sigma_{N^{(d)}}}(\infty)| \geq \sqrt{\mu_d(M)}\}| \leq 3\sqrt{\mu_d(M)}M'.$$

By the pigeonhole principle, we can thus find an interval  $[M'', M'' + 1/10\sqrt{\mu_d(M)}]$  in  $[0, M')$  such that

$$|\sigma_{N^{(d)}}(Q_n) - \overline{\sigma_{N^{(d)}}}(\infty)| < \sqrt{\mu_d(M)} \text{ for all } n \in [M'', M'' + 1/10\sqrt{\mu_d(M)}].$$

Since  $\overline{\sigma_{N^{(d)}}}(\infty) \geq 1/2$  and  $\mu_d(M)$  is small, we see in particular that  $Q_n$  contains at least one saturated component and in particular is not black.

The number of possible colors of  $Q_n$  is at most  $2^{N_1 \dots N_d M}$ . Thus, if  $c$  is sufficiently small (and hence  $1/10\sqrt{\mu_d(M)}$  sufficiently large) depending on  $k$  and  $2^{N_1 \dots N_d M}$ , we see from van der Waerden's theorem that we can find an increasing progression  $n, \dots, n + (k-1)r$  inside  $[M'', M'' + 1/10\sqrt{\mu_d(M)}]$  such that all the progressions  $Q_n, \dots, Q_{n+(k-1)r}$  have the same color, which is not black by the preceding discussion. In particular inside  $Q_n \cup \dots \cup Q_{n+(k-1)r}$  we can find an increasing arithmetic progression of length  $k$ , and we are done.  $k$  as desired.  $\blacksquare$

We will assume that the asymptotic upper saturation is quite large, say

$$\overline{\sigma_{N^{(d)}}}(\infty) \geq 1 - \varepsilon_d \geq 1 - 1/100k \quad (2)$$

for some small  $0 \leq \varepsilon_d < 1/100k$ . This means that we can find arbitrarily large rank  $d+1$  increasing progressions which are almost entirely saturated (up to an exceptional set of density at most  $\varepsilon_d$ ). In particular we see that  $\Sigma_{N^{(d)}}$  is non-empty.

Given any saturated color  $c \in \Sigma_{N^{(d)}}$ , we may define the *color density*

$$\delta_c(P) := \frac{1}{N_{d+1}} |\{0 \leq j < N_{d+1} : c(P_j) = c\}|$$

of any increasing progression  $P$  of dimensions  $(N_1, \dots, N_{d+1})$  (again note that this depends only on the color  $c(P)$  of  $P$ , and not on  $P$  itself), as well as the upper color density numbers

$$\overline{\delta}_c(N_{d+1}) := \sup_{P: |\sigma_{N^{(d)}}(P) - \overline{\sigma_{N^{(d)}}}(\infty)| \leq \sqrt{\mu_d(N_{d+1})}} \delta_c(P). \quad (3)$$

Note that we have restricted the supremum to progressions  $P$  which are quite saturated; this set is nonempty by definition of  $\sigma_{N^{(d)}}(N_{d+1})$  and  $\mu_d(N_{d+1})$ . We now define the asymptotic upper color density

$$\overline{\delta}_c(\infty) := \limsup_{N_{d+1} \rightarrow \infty} \overline{\delta}_c(N_{d+1}). \quad (4)$$

Since

$$\overline{\sigma_{N^{(d)}}}(\infty) \leq \sum_{c \in \Sigma_{N^{(d)}}} \overline{\delta_c}(\infty)$$

we see from the pigeonhole principle and (2) that there exists a color  $c_d \in \Sigma_{N^{(d)}}$  (depending of course on  $N^{(d)}$  and  $\Sigma_{N^{(d)}}$ ) such that

$$\overline{\delta_{c_d}}(\infty) \geq \frac{1 - \varepsilon_d}{|\Sigma_{N^{(d)}}|} \geq \frac{1}{2|2^{[0, N^{(d)}]}|}. \quad (5)$$

Let us fix this color  $c_d$ , and refer to it as a *perfect* color (the notation is essentially from [3]); one can think of this color as the most “popular” of the saturated colors. In practice, we will be able to ensure that this color is not black (either by choosing it specifically (in the  $d = 0$  case), or by assuming inductively that all the saturated colors are non-black). The condition (5) ensures that the perfect color is attained quite often; while the denominator  $|2^{[0, N^{(d)}]}|$  on the right-hand side of (5) looks quite large, we shall eventually choose  $N_{d+1}$  (and related quantities) much larger than this, to the extent that this denominator ends up being negligible.

Let us now pick a small  $0 < \varepsilon_{d+1} < 1/100k$  (which will in practice be much smaller than  $\varepsilon_d$  or  $\overline{\delta_{c_d}}(\infty)$ ). By (1), (4), we can find arbitrarily large  $N_{d+1}$  such that

$$|\overline{\delta_{c_d}}(N_{d+1}) - \overline{\delta_{c_d}}(\infty)| \leq \varepsilon_{d+1}^8 \quad (6)$$

and

$$\mu_d(N_{d+1}) \leq \varepsilon_{d+1}^8. \quad (7)$$

Suppose we choose such a  $N_{d+1}$ . We then form the dimensions  $N^{(d+1)} := (N_1, \dots, N_d, N_{d+1})$ . If  $P$  is an increasing progression of dimensions  $N^{(d+1)}$ , we say that  $c(P)$  is *saturated* if

$$|\sigma_{N^{(d)}}(P) - \overline{\sigma_{N^{(d)}}}(\infty)| \leq \sqrt{\mu_d(N_{d+1})} \quad (8)$$

and

$$|\delta_{c_d}(P) - \overline{\delta_{c_d}}(\infty)| \leq \varepsilon_{d+1}. \quad (9)$$

Note that this definition is well-defined as the quantities  $\sigma_{N^{(d)}}(P)$  and  $\delta_{c_d}(P)$  depend only on the color  $c(P)$  of  $P$  and not of  $P$  itself. Informally, a progression of rank  $d + 1$  is saturated if almost all of its components are saturated, and as many of its components as possible have the perfect color. We can thus define a set  $\Sigma_{N^{(d+1)}} \subseteq 2^{[0, N^{(d+1)})}$  of saturated colors, which in turn induces an asymptotic upper saturation  $\overline{\sigma_{N^{(d+1)}}}$  as before. We observe also that (8), (2), (7) implies that

$$\sigma_{N^{(d)}}(P) \geq 1 - 2\varepsilon_{d+1} \quad (10)$$

(for instance).

**Lemma 3.3.** *We have*

$$\overline{\sigma_{N^{(d+1)}}}(\infty) \geq 1 - \varepsilon_{d+1}.$$

(In other words, (2) holds for rank  $d + 1$ ).

**Proof** Let  $\kappa > 0$  and  $M > 0$  be arbitrary. It will suffice to show that there exists an  $N_{d+2} \geq M$  and an increasing progression  $P$  of rank  $d+2$  and dimensions  $(N_1, \dots, N_{d+2})$  such that

$$\sigma_{N^{(d+1)}}(P) \geq 1 - \varepsilon_{d+1}. \quad (11)$$

By (1), (4), we can find  $M' \geq N_{d+1} \max(M, \frac{1}{\mu_d(N_{d+1})^2})$  so large that

$$\overline{\sigma_{N^{(d)}}}(M') = \overline{\sigma_{N^{(d)}}}(\infty) + O(\mu_d(N_{d+1})^2) \quad (12)$$

and then an increasing progression  $Q$  of rank  $d+1$  and dimensions  $(N_1, \dots, N_d, M')$  such that

$$\delta_{c_d}(Q) = \overline{\delta_{c_d}}(\infty) + O(\mu_d(N_{d+1})) \quad (13)$$

and

$$\sigma_{N^{(d)}}(Q) = \overline{\sigma_{N^{(d)}}}(\infty) + O(\mu_d(N_{d+1})).$$

By shrinking  $Q$  a little bit we may assume that  $M'$  is a multiple of  $N_{d+1}$  (the error incurred by this is at most  $O(N_{d+1}/M') = O(\mu_d(N_{d+1})\varepsilon_{d+1})$ ), say  $M' = N_{d+1}N_{d+2}$ . We can then view  $Q$  not as an increasing progression of rank  $d+1$  and dimensions  $(N_1, \dots, N_d, M')$ , but rather as an increasing progression (which we shall call  $P$  to distinguish it from  $Q$ , even though as sets of integers they are identical) of rank  $d+2$  and dimensions  $(N_1, \dots, N_d, N_{d+1}, N_{d+2})$ . Since

$$\sigma_{N^{(d)}}(Q) = \frac{1}{N_{d+2}} \sum_{n_{d+2} \in [0, N_{d+2})} \sigma_{N^{(d)}}(P_{n_{d+2}})$$

we thus have

$$\frac{1}{N_{d+2}} \sum_{n_{d+2} \in [0, N_{d+2})} (\sigma_{N^{(d)}}(P_{n_{d+2}}) - \overline{\sigma_{N^{(d)}}}(\infty)) = O(\mu_d(N_{d+1})).$$

But from definition of  $\mu_d(N_{d+1})$  we have

$$\frac{1}{N_{d+2}} \sum_{n_{d+2} \in [0, N_{d+2})} (\sigma_{N^{(d)}}(P_{n_{d+2}}) - \overline{\sigma_{N^{(d)}}}(\infty))_+ \leq \mu_d(N_{d+1})$$

and hence

$$\frac{1}{N_{d+2}} \sum_{n_{d+2} \in [0, N_{d+2})} |\sigma_{N^{(d)}}(P_{n_{d+2}}) - \overline{\sigma_{N^{(d)}}}(\infty)| = O(\mu_d(N_{d+1})).$$

By Chebyshev's inequality and (7) we thus have

$$\{n_{d+2} \in [0, N_{d+2}) : |\sigma_{N^{(d)}}(P_{n_{d+2}}) - \overline{\sigma_{N^{(d)}}}(\infty)| \geq \sqrt{\mu_d(N_{d+1})}\} \leq O(\varepsilon_{d+1}^4)N_{d+2}. \quad (14)$$

Now we argue using  $\delta_{c_d}$  instead of  $\sigma_{N^{(d)}}$ . From (13), (7) we have

$$\frac{1}{N_{d+2}} \sum_{n_{d+2} \in [0, N_{d+2})} (\delta_{c_d}(P_{n_{d+2}}) - \delta_{c_d}(\infty)) = O(\varepsilon_{d+1}^8),$$

so by (14)

$$\frac{1}{N_{d+2}} \sum_{n_{d+2} \in [0, N_{d+2}] : |\sigma_{N^{(d)}}(P_{n_{d+2}}) - \sigma_{N^{(d)}}(\infty)| < \sqrt{\mu_d(N_{d+1})}} (\delta_{c_d}(P_{n_{d+2}}) - \delta_{c_d}(\infty)) = O(\varepsilon_{d+1}^4).$$

But by (6) we have

$$\frac{1}{N_{d+2}} \sum_{n_{d+2} \in [0, N_{d+2}] : |\sigma_{N^{(d)}}(P_{n_{d+2}}) - \sigma_{N^{(d)}}(\infty)| < \sqrt{\mu_d(N_{d+1})}} (\delta_{c_d}(P_{n_{d+2}}) - \delta_{c_d}(\infty))_+ = O(\varepsilon_{d+1}^8)$$

and so

$$\frac{1}{N_{d+2}} \sum_{n_{d+2} \in [0, N_{d+2}] : |\sigma_{N^{(d)}}(P_{n_{d+2}}) - \sigma_{N^{(d)}}(\infty)| < \sqrt{\mu_d(N_{d+1})}} |\delta_{c_d}(P_{n_{d+2}}) - \delta_{c_d}(\infty)| = O(\varepsilon_{d+1}^4)$$

and so by Chebyshev

$$|\{n_{d+2} \in [0, N_{d+2}] : |\sigma_{N^{(d)}}(P_{n_{d+2}}) - \sigma_{N^{(d)}}(\infty)| < \sqrt{\mu_d(N_{d+1})}; |\delta_{c_d}(P_{n_{d+2}}) - \delta_{c_d}(\infty)| \geq \varepsilon_{d+1}\}| = O(\varepsilon_{d+1}^3)N_{d+2}.$$

Combining this with (14) and the definition of saturation we see that

$$|\{n_{d+2} \in [0, N_{d+2}] : P_{n_{d+2}} \text{ is not saturated}\}| = O(\varepsilon_{d+1}^3)N_{d+2}$$

and hence (11) follows.  $\blacksquare$

This lemma will allow us to iteratively construct a sequence  $N_1 \ll N_2 \ll \dots$  of integers, with the associated dimensions  $N^{(d)} := (N_1, \dots, N_d)$ , a set  $\Sigma_{N^{(d)}} \subset 2^{[0, N^{(d)})}$  of saturated colors, and a collection of perfect colors  $c_d \in \Sigma_{N^{(d)}}$  for each  $d$ ; we shall formalize this construction in the next section.

#### 4. SETTING UP THE DIMENSIONS

We shall need a *growth function*  $F_0 : \mathbf{R}^+ \rightarrow \mathbf{R}^+$ , depending on  $k$  and  $\delta$  to be chosen later. One should think of this function as a rapidly increasing function, such as  $F_0(n) := 2^{2^{kn/\delta}}$ , though in practice we shall need a much faster function than this (it has to grow faster than the bounds one obtains from both the Szemerédi regularity lemma and the van der Waerden theorem). It will be important that our bounds will not depend on this extremely rapidly growing function. We will however make the (very mild) assumption that  $F_0(n)$  grows faster than  $n$ , e.g.  $F_0(n) \geq 2^n$ .

As mentioned in the previous section, we wish to construct a sequence  $N_1 \ll N_2 \ll \dots$  of integers as well as appropriate notions of saturated color and perfect color for each  $d \geq 0$  (actually, at the end of the day we will only need this for  $d \leq 2^k$ ). We will also obtain error estimates of the form (2) for some  $0 < \varepsilon_d < 1/100k$ , as well as lower bounds on the asymptotic perfect color density.



**4.1. The  $d = 0$  constructions.** We begin with the  $d = 0$  case. Here  $N^{(0)} = ()$  is the empty tuple, and progressions of rank 0 are just singleton sets of integers. Such singletons are either *black* (if they are not contained in  $A$ ) or *white* (if they are contained in  $A$ ), giving two colors. We declare both of them to be saturated (thus  $\Sigma_{N^{(0)}} := 2^{[0, N^{(0)})} = \{\text{black}, \text{white}\}$ ) and we declare  $c_0 := \text{white}$  to be the perfect color.

Since every color is saturated, we have  $\overline{\sigma_{N^{(0)}}}(\infty) = 1$  and so if we set  $\varepsilon_0 := 0$  then (2) is satisfied.

$$\delta_{N^{(0)}, c_0}([-N, N]) = \frac{|A \cap [-N, N]|}{|[-N, N]|}$$

and  $A$  has upper density at least  $\delta$ , we see that

$$\delta_{c_0} \geq \delta. \quad (15)$$

**4.2. The general case.** Now we assume inductively that the quantities  $N^{(d)}, \varepsilon_d, \Sigma_{N^{(d)}}, c_d$  have already been constructed for some  $d \geq 0$ , and turn now to the construction for  $d + 1$ .

Using (1), (4), we can find a length

$$L_{d+1} \geq F_0(N_1 \dots N_d) \quad (16)$$

such that

$$|\overline{\delta_{c_d}}(L_{d+1}) - \overline{\delta_{c_d}}(\infty)| \leq 1/F_0(N_1 \dots N_d) \quad (17)$$

and

$$|\mu_d(L_{d+1})| \leq 1/F_0(N_1 \dots N_d)^2. \quad (18)$$

We then define a small number  $0 < \varepsilon_{d+1} \ll 1$  defined as

$$\varepsilon_{d+1} := (100kF_0(L_{d+1}))^{-100k} \quad (19)$$

(so in particular  $\varepsilon_{d+1} < 1/100k$ ), and then using (4) again we can find a length

$$N_{d+1} \geq F_0(1/\varepsilon_{d+1}) \quad (20)$$

such that (6), (7) hold. We then use this length  $N_{d+1}$  to construct the dimensions  $N^{(d+1)} := (N_1, \dots, N_{d+1})$ , and then define the saturated color set  $\Sigma_{N^{(d+1)}}$  and the perfect color  $c_{d+1}$  as in the previous section. We can then iterate this construction indefinitely to create  $N^{(d)}, \varepsilon_d, \Sigma_{N^{(d)}}, c_d$  for all  $d$ .

Note that the quantities  $L_{d+1}, 1/\varepsilon_{d+1}, N_{d+1}$  are widely separated in magnitude, indeed we have

$$\dots \ll N_d \ll L_{d+1} \ll 1/\varepsilon_{d+1} \ll N_{d+1} \ll L_{d+2} \ll \dots$$

where  $A \ll B$  denotes the statement that  $B \geq F_0(A)$ . The fact that  $F_0$  is extremely rapidly growing will mean that any reasonable quantity involving variables to the left of this hierarchy will be dominated by any reasonable quantity involving variables to the right of this hierarchy. The parameter  $L_{d+1}$  is an intermediate dimension which is huge compared to  $N_d$  but miniscule compared to  $N_{d+1}$ ; the point is that when we are working in a progression of dimensions  $N^{(d+1)}$  we shall be easily able to introduce a free parameter ranging over an interval of length  $L_{d+1}$

(or even  $F_0(L_{d+1})$ , which would be pretty useless if we stayed in rank  $d + 1$ , but if we then drop rank to  $d$ , this additional parameter becomes extremely useful - roughly speaking, it (combined with the Szemerédi regularity lemma and the van der Waerden theorem, and the fact that  $L_{d+1}$  has upper perfect color density close to the asymptotic limit) can be used to upgrade the existence of progressions of length  $i$  to progressions of length  $i + 1$  (although this “uses” up our parameter and we have to go create it again). Thus we shall be able to “spend” the large rank in our progressions, one rank at a time, to eventually obtain progressions of length  $k$ . This strategy may appear extremely convoluted, but it in fact appears to be the minimally complex strategy that would combine all these ingredients to yield a successful proof (and avoid all the difficulties associated with the bounds for the regularity lemma and van der Waerden theorem being incredibly poor).

We now observe that none of the perfect colors  $c_d$  are black. For  $d = 0$  this is clear from construction. For  $d = 1$ , we see from (19) that

$$\varepsilon_1 < \delta/2$$

if  $F_0$  is chosen sufficiently large depending on  $\delta$ , and so from (9) we see that  $\delta_{c_0}(P) > 0$  whenever  $P$  is a rank 1 saturated increasing progression. This implies that all the saturated colors in  $\Sigma_{N^{(1)}}$  are non-black, and in particular the perfect color  $c_1$  is not black.

Now suppose inductively that  $c_d$  is known to be non-black for some  $d \geq 1$ . From (19) we certainly have

$$\varepsilon_{d+1} < \frac{1}{4|2^{[0, N^{(d)}]}|}$$

if  $F_0$  is chosen sufficiently large, and so by (5), (9) we see that  $\delta_{c_d}(P) > 0$  whenever  $P$  is a rank  $d + 1$  saturated increasing progression, thus as before we see that all the saturated colors in  $\Sigma_{N^{(d+1)}}$  are non-black, and in particular the perfect color  $c_{d+1}$  is also non-black.

Since none of the saturated colors in  $\Sigma_{N^{(d)}}$  are black for  $d \geq 1$ , we see from Lemma 3.2 that we may assume

$$\mu_d(L_{d+1}) > c(k, N_1, \dots, N_d, L_{d+1}). \quad (21)$$

Note from the above construction that we have in fact shown

$$\varepsilon_{d+1} < \frac{1}{2}\delta_{c_d}(\infty)$$

and hence

$$\frac{1}{2}\delta_{c_d}(\infty) \leq \delta_{c_d}(P) \leq 2\delta_{c_d}(\infty) \quad (22)$$

whenever  $P$  is a rank  $d + 1$  saturated increasing progression.

If  $d \geq 0$ ,  $0 \leq i \leq k$ , and  $P$  is an increasing progression of dimensions  $(N_1, \dots, N_d, k + 1)$ , we say that  $P$  is *perfect of order  $i$*  if the first  $i$  components  $P_0, \dots, P_{i-1}$  all have the perfect color  $c_d$ . Clearly all increasing progressions of dimension  $(N_1, \dots, N_d, k + 1)$  are perfect of order 0; in order to prove Theorem 1.1 we will seek to find progressions which are perfect of order  $k$  and have dimension  $k + 1$

(actually we can eventually obtain progressions which are perfect of order  $k$  and have dimension  $(N_1, \dots, N_d, k+1)$  for any  $k$ ).

We mentioned earlier that our basic strategy is to start with a large rank progression, and spend those ranks in order to improve the order  $i$  of perfection. We now turn to setting up the notation required to execute this strategy.

## 5. HOMOGENEOUS PROGRESSIONS, AND WELL-ARRANGED SEQUENCES OF PROGRESSIONS

Let  $d \geq 1$ , and let  $P$  be a rank  $d+1$  increasing progression of dimensions  $(N_1, \dots, N_d, k+1)$ . Thus  $P$  consists of  $k+1$  components  $P_0, \dots, P_k$  which are rank  $d$  increasing progressions of dimensions  $(N_1, \dots, N_d)$ , and each component  $P_j$  in turn consists of  $N_d$  components  $P_{j,0}, \dots, P_{j,N_d-1}$  which are rank  $d-1$  increasing progressions of dimensions  $(N_1, \dots, N_{d-1})$ .

Let us say that  $P$  is *completely saturated* if each of its  $k+1$  components  $P_0, \dots, P_k$  are saturated. By (10) this means that

$$\frac{1}{N_d} |\{0 \leq n_d < N_d : P_{j,n_d} \text{ is saturated}\}| \geq 1 - 2\varepsilon_{d-1} \text{ for all } 0 \leq j \leq k$$

and from (22) we similarly have

$$\frac{1}{N_d} |\{0 \leq n_d < N_d : P_{j,n_d} \text{ has the perfect color } c_{d-1}\}| = \delta_{c_{d-1}}(\infty) + O(\varepsilon_d) \text{ for all } 0 \leq j \leq k. \quad (23)$$

If  $m_d, m_d+r, \dots, m_d+(k-1)r$  is a progression of integers in  $[0, N_d)$  (not necessarily increasing or proper), we can define the associated *subprogression*  $P_{(m_d,r)} \subset P$  of  $P$  to be the set

$$P_{(m_d,r)} = P_{0,m_d} \cup P_{1,m_d+r} \cup \dots \cup P_{k-1,m_d+kr};$$

observe that this is a rank  $d$  increasing progression of dimensions  $(N_1, \dots, N_{d-1}, k+1)$ .

We shall need the counting functions

$$f_{P,i,j}(n_d) := |\{P_{(m_d,r)} \subset P : P_{(m_d,r)} \text{ is perfect of order } i \text{ and } m_d + jr = n_d\}|$$

defined for all  $0 \leq i \leq j \leq k$  and  $n_d \in [0, N_d)$ ; thus  $f_{P,i,j}(n_d)$  counts the number of subprogressions  $P_{(m_d,r)}$  which are perfect of order  $i$  and contain the specific component  $P_{j,n_d}$ . Thus, for instance, the existence of a subprogression which is perfect of order  $i$  is equivalent to  $f_{P,i,j}(n_d)$  being non-zero for some  $j$  and some  $n_d$ . This quantity is clearly decreasing in  $i$ ,  $f_{P,i+1,j}(n_d) \leq f_{P,i,j}(n_d)$ , and is easy to compute when  $i = 0$ :

**Lemma 5.1.** *If  $n_d \in [0, N_d)$  and  $0 \leq j \leq k$ , then we have the upper bound*

$$f_{P,0,j}(n_d) \leq N_d. \quad (24)$$

If in addition  $n_d$  lies in the middle third  $[N_d/3, 2N_d/3)$  of  $[0, N_d)$ , then we can complement this upper bound with the lower bound

$$f_{P,0,j}(n_d) \frac{1}{10k} N_d.$$

**Proof** The estimate (24) just reflects the fact that once  $n_d$  is fixed, one can specify the progression  $P_{(m_d,r)}$  by fixing  $m_d + j'r$  for some  $j' \neq j$ , but this ranges in  $[0, N_d)$ . The lower bound follows for instance by observing that any  $|r| \leq 1/20k$  will generate a progression  $P_{(m_d,r)} \subset P$  with  $m_d + jr = n_d$ , which is then automatically perfect of order 0.  $\blacksquare$

Recall from (23) that the proportion of components  $P_{j,n_d}$  which have the perfect color is roughly  $\delta_{c_{d-1}}(\infty)$ , or in other words the probability that a randomly chosen component is perfectly colored is roughly  $\delta_{c_{d-1}}(\infty)$ . If the perfectly colored components were distributed “independently of each other”, then probabilistic heuristics would then suggest that  $f_{P,i,j}(n_d) \approx \delta_{c_{d-1}}^i(\infty) f_{P,0,j}(n_d)$ , and so in particular for  $n_d$  in the middle third  $[N_d/3, 2N_d/3)$  one would expect  $f_{P,i,j}(n_d)$  to be roughly of the order of  $\delta_{c_{d-1}}^i(\infty) N_d$  for most  $n_d$ . Now, *a priori* there is no reason why one should expect this property (which is broadly analogous to the concepts of “weak mixing” used in the Furstenberg proof of Szemerédi’s theorem, “Gowers uniformity” in the Gowers’ proof, or “ $\varepsilon$ -regularity” in the proof of the Szemerédi regularity lemma) should hold for a typical progression  $P$ . If however we have a long sequence of such progressions  $P$  which are arranged in a special way (we will make this concept more precise later), it will turn out that this type of property will in fact hold for at least one progression in a sequence, by means of the Szemerédi regularity lemma combined with van der Waerden’s theorem. This observation is crucial in letting us induct on the parameter  $i$  and create increasingly perfect progressions.

We must first formalize the notion of what it means for  $f_{P,i,j}(n_d)$  to be roughly of the order of  $\delta_{c_{d-1}}^i(\infty) N_d$  for most  $n_d$ . This will be done using the concept of a *homogeneous* progression (again, the notation is from [3]).

**Definition 5.2.** Let  $d \geq 1$ , let  $0 \leq i \leq k$ , and let  $P$  be a rank  $d + 1$  increasing progression of dimensions  $(N_1, \dots, N_d, k + 1)$ . We say that  $P$  is *homogeneous of order  $i$*  if it obeys the following two properties:

- $P$  is completely saturated (i.e. each of the  $k + 1$  components  $P_0, \dots, P_k$  is saturated); and
- For all  $i \leq j \leq k$ , we have

$$|\{n_d \in [0, N_d) : f_{P,i,j}(n_d) > (100k)^i \delta_{c_{d-1}}(\infty)^i N_d\}| \leq i \varepsilon_{d-1} \delta_{c_{d-1}}(\infty)^{k+1} N_d \quad (25)$$

and

$$|\{n_d \in [N_d/3, 2N_d/3) : f_{P,i,j}(n_d) < (100k)^{-i-1} \delta_{c_{d-1}}(\infty)^i N_d\}| \leq i \varepsilon_{d-1} \delta_{c_{d-1}}(\infty)^{k+1} N_d. \quad (26)$$

Thus for almost all  $n_d$  in  $[0, N_d)$ , and especially in the middle third  $[N_d/3, 2N_d/3)$ , the counting function  $f_{P,i,j}(n_d)$  is fairly close to its expected value of

$\delta_{c_{d-1}(\infty)}^i N_d$  (the factors of  $(100k)^{\pm i}$  are extremely negligible compared to the very huge numbers and very small being manipulated here).

Thus for instance, any progression  $P$  of dimensions  $(N_1, \dots, N_d, k+1)$  which is completely saturated will automatically be homogeneous of order 0, thanks to Lemma 5.1. It will be of interest to obtain progressions which are homogeneous of as high an order  $i$  as possible, as this will imply (by (26)) that there are many subprogressions  $P_{(m_d, r)}$  which are perfect of order  $i$  (this is the place where we drop a rank). Furthermore, we will later show (by means of the van der Waerden theorem and the Szemerédi regularity lemma) that given a sufficiently long sequence of progressions  $P$  which are homogeneous of order  $i$  and are arranged in a special way, one of them will be homogeneous of order  $i+1$ . The combination of these two facts (by a simple induction argument) will eventually give us progressions which are perfect of any order  $0 \leq i \leq k$ , which in particular will give Proposition ??.

To make this argument rigorous we need to pin down what we mean by a sequence of progressions  $P$  being “arranged in a special way”. We shall do this via the following somewhat technical definition (which is forced on us by the inductive argument we shall use).

**Definition 5.3.** Let  $d \geq 1$ , let  $\Omega \subset [0, k]$ . We define a *well-arranged sequence of progressions*  $P[]$  of rank  $d+1$  and *parameterization*  $\Omega$  to be an assignment  $P[(\lambda_j)_{j \in \Omega}]$  of a rank  $d+1$  increasing progression of dimensions  $(N_1, \dots, N_d, k+1)$  to each  $\Omega$ -tuple  $(\lambda_j)_{j \in \Omega} \in [0, F_0(L_d)]^\Omega$  (which we refer to as the *parameters* of the sequence) such that

- (i) (Complete saturation property) For each  $(\lambda_j)_{j \in \Omega} \in [0, L_d]^\Omega$ , the progression  $P[(\lambda_j)_{j \in \Omega}]$  is completely saturated.
- (ii) (Independence property) For any  $0 \leq i' \leq k$ , the color  $c(P[(\lambda_j)_{j \in \Omega}]_{i'}) \in 2^{\{0, N^{(d)}\}}$  of the  $i'$  component of  $P[(\lambda_j)_{j \in \Omega}]$  is determined entirely by those  $\lambda_j$  for which  $j \leq i'$ ; to put this another way, the parameter  $\lambda_j$  influences the colors of the  $j, \dots, k$  components of  $P[(\lambda_j)_{j \in \Omega}]$  but has no effect on the colors of the  $0, \dots, j-1$  components.
- (iii) (Arithmetic structure) For each  $j_0 \in \Omega$ , if we fix all of the  $\lambda_j$  indices except for the  $\lambda_{j_0}$  index, then the components  $P[(\lambda_j)_{j \in \Omega}]_{j_0}$ , as  $\lambda_{j_0}$  ranges from 0 to  $L_d - 1$ , form an increasing rank  $d$  progression of dimensions  $(N_1, \dots, N_{d-1}, L_d)$ .

Note that in the case  $\Omega = \emptyset$ , then there are no free parameters, and any completely saturated rank  $d+1$  progression  $P$  of dimensions  $(N_1, \dots, N_d, k+1)$  is automatically well-arranged. The independence assumptions (ii) are necessary in order for us to manipulate the parameters  $\lambda_j, j \in \Omega$  more or less independently; the arithmetic structure assumptions (iii) are needed in order for us to apply van der Waerden's theorem at a key juncture.

The proof of Proposition ?? centers around the following two key propositions, which correspond roughly to Lemma 6 and Lemma 5 from [3].

**Proposition 5.4.** [3, Lemma 6] *Let  $d \geq 2$ ,  $0 \leq i \leq k$ , and let  $\Omega \subset [i, k]$  be such that  $i \notin \Omega$  (i.e.  $\Omega \subset [i+1, k]$ ). Suppose that we have a well-arranged sequence of progressions  $P[]$  of rank  $d+1$  and parameterization  $\Omega$  such that every progression  $P[(\lambda_j)_{j \in \Omega}]$  in this sequence is homogeneous of order  $i$ . Then there exists another well-arranged sequence of progressions  $Q[]$  of rank  $d$  and parameterization  $\Omega \cup \{i\}$ .*

**Proposition 5.5.** [3, Lemma 5] *Let  $d \geq 1$ ,  $0 \leq i \leq k$ , and let  $\Omega \subset [i, k]$  be such that  $i \in \Omega$ . Suppose that we have a well-arranged sequence of progressions  $P[]$  of rank  $d+1$  and parameterization  $\Omega$  such that every progression  $P[(\lambda_j)_{j \in \Omega}]$  in this sequence is homogeneous of order  $i$ . Then there exists another well-arranged sequence of progressions  $Q[]$  of rank  $d+1$  and parameterization  $\Omega \setminus \{i\}$ , such that every progression  $Q[(\lambda_j)_{j \in \Omega \setminus \{i\}}]$  in this sequence is homogeneous of order  $i+1$ .*

Proposition 5.4 allows us to convert homogeneity of order  $i$  to an additional parameterization, by the  $\lambda_i$  parameter, at the cost of dropping one rank; Proposition 5.5 allows us to spend such a  $\lambda_i$  parameter to upgrade homogeneity of order  $i$  to homogeneity of order  $i+1$ , keeping the rank constant. These two propositions will be proven in the next two sections. Assuming them for the moment, the two propositions combine to form the following key fact, which is essentially Fact 12 from [3]:

**Corollary 5.6.** [3, Fact 12] *Let  $0 \leq i \leq k$ ,  $d \geq 2^i$ , and let  $\Omega \subset [i, k]$ . Suppose that we have a well-arranged sequence of progressions  $P[]$  of rank  $d+1$  and parameterization  $\Omega$ . Then there exists another well-arranged sequence of progression  $Q[]$  of rank  $d - 2^i + 2$  and parameterization  $\Omega$  such that every progression  $Q[(\lambda_j)_{j \in \Omega}]$  in this sequence is homogeneous of order  $i$ .*

**Proof** We prove this by induction on  $i$ . The claim is trivial for  $i = 0$  since, by property (i) of a well-arranged sequence, each  $P[(\lambda_j)_{j \in \Omega}]$  is completely saturated and hence homogeneous of order 0, so we may simply take  $Q := P$ .

Now suppose that  $1 \leq i \leq k$ , and the claim has already been proven for  $i-1$ . Let  $P[]$  be a well-arranged sequence of progressions  $P[]$  of rank  $d+1$  and parameterization  $\Omega$  for some  $\Omega \subset [i, k]$ . By the induction hypothesis, we can find a well-arranged sequence of progressions  $P'[]$  of rank  $d - 2^{i-1} + 2$  and parameterization  $\Omega$  such that each  $P'[(\lambda_j)_{j \in \Omega}]$  in this sequence is homogeneous of order  $i-1$ . If we then apply Proposition 5.4 to this sequence, we obtain another well-arranged sequence of progressions  $P''[]$  of rank  $d - 2^{i-1} + 1$  and parameterization  $\Omega \cup \{i-1\}$ . Applying the induction hypothesis again, we can then find another well-arranged sequence of progressions  $P'''[]$  of rank  $d - 2^i + 2$  and parameterization  $\Omega \cup \{i-1\}$ , such that each  $P'''[(\lambda_j)_{j \in \Omega}]$  in this sequence is homogeneous of order  $i-1$ . Applying Proposition 5.5, we finally obtain a well-arranged sequence of progressions  $Q[]$  of rank  $d - 2^i + 2$  and parameterization  $\Omega$ , such that every progression  $Q[(\lambda_j)_{j \in \Omega}]$  in this sequence is homogeneous of order  $i$ , thus closing the induction.  $\blacksquare$

We can now prove Theorem 1.1, as follows.

**Proof** [Proof of Theorem 1.1] Let  $d := 2^k$ . From (2) we know that

$$\overline{\sigma_{N^{(d)}}(\infty)} \geq 1 - 1/100k$$

and hence there exists some  $L \geq 100k$  and some increasing progression  $P$  of rank  $d + 1$  and dimensions  $(N_1, \dots, N_d, L)$  such that

$$\sigma_{N^{(d)}}(P) \geq 1 - 1/100k,$$

i.e. all the components  $P_0, \dots, P_{L-1}$  of  $P$  are saturated with at most  $L/100k$  exceptions. Applying Lemma 1.2, we can thus find an increasing subprogression  $P' = P_l \cup P_{l+r}, \dots, P_{l+kr}$  which is completely saturated. This progression  $P'$  has rank  $d + 1$  and dimensions  $(N_1, \dots, N_d, k + 1)$ ; since it is completely saturated, it can be viewed as a well-arranged sequence  $P' = P'[]$  of rank  $d + 1$  parameterized by the empty set  $\Omega = \emptyset$ . Applying Corollary 5.6 with  $i = k$  and  $\Omega = \emptyset$  we can thus find a completely saturated increasing progression  $Q'$  of rank 2 and dimensions  $(N_1, k + 1)$  which is homogeneous of degree  $k$ . From (26) we thus see that there exists an  $n_1 \in [N_1/3, 2N_1/3)$  such that  $f_{Q', k, k}(n_1) > 0$ , which implies in particular that there exists an increasing progression  $R$  of length  $k + 1$  which is perfect of order  $k$ , i.e. the first  $k$  elements of this progression lie in  $A$ , as desired. ■

The only thing left to do is prove Propositions 5.4 and 5.5.

## 6. PROOF OF PROPOSITION 5.4

We now prove Proposition 5.4, which is the more elementary of the two propositions (it does not require either the van der Waerden theorem or the Szemerédi regularity lemma); the principal difficulty is to ensure that all the relevant components of the various progressions involved are saturated.

Fix  $P, i, \Omega$ . If  $m_d, m_d + r, \dots, m_d + kr$  is a progression of integers in  $[0, N_d)$ , let us call the pair  $(m_d, r)$  *perfect of order  $i$*  if the subprogression  $P[(\lambda_j)_{j \in \Omega}]_{(m_d, r)}$  is perfect of order  $i$  for some  $(\lambda_j)_{j \in \Omega} \in [0, F_0(L_{d-1}))^\Omega$ ; we have truncated the range of the parameters  $\lambda_j$  from  $[0, F_0(L_d))$  to the much smaller (and hence more tractable) range of  $[0, F_0(L_{d-1}))$  as we shall eventually drop rank by 1. Actually, the choice of the  $\lambda_j$  is irrelevant since the concept of perfection will only depend on the color of first  $i$  components  $P[(\lambda_j)_{j \in \Omega}]_{0, m_d}, \dots, P[(\lambda_j)_{j \in \Omega}]_{i-1, m_d + (i-1)r}$ , whereas by the hypothesis  $\Omega \subset [i, k]$  and by the independence property (ii) of a well-arranged sequence, the colors of the first  $i$  components do not depend on any of the parameters  $\lambda_j$  for  $j \in \Omega$ . Let us call the pair  $(m_d, r)$  *totally saturated* if the subprogression  $P[(\lambda_j)_{j \in \Omega}]_{(m_d, r)}$  is completely saturated for every  $(\lambda_j)_{j \in \Omega} \in [0, F_0(L_{d-1}))^\Omega$ . Unlike the concept of perfection of order  $i$ , the parameters  $\lambda_j$  do play a non-trivial role in this concept, as the concept of saturation involves the colors of all  $k + 1$  components of  $P[(\lambda_j)_{j \in \Omega}]_{(m_d, r)}$ . However, since all the  $P[(\lambda_j)_{j \in \Omega}]$  are completely saturated (by property (i) of a well-arranged sequence), the estimate (10) asserts that almost all components  $P[(\lambda_j)_{j \in \Omega}]_{i', n_d}$  will be saturated:

$$\frac{1}{N_d} |\{n_d \in [0, N_d) : P[(\lambda_j)_{j \in \Omega}]_{i', n_d} \text{ is not saturated}\}| \leq 2\varepsilon_{d-1} \text{ for all } (\lambda_j)_{j \in \Omega} \in [0, F_0(L_{d-1}))^\Omega, 0 \leq i' \leq k. \quad (27)$$

This basically means that the property of being completely saturated is easy to attain. This, combined with (26) makes the following lemma fairly reasonable:

**Lemma 6.1.** *Call an element  $n_d \in [N_d/3, 2N_d/3]$  good if there exists a progression  $m_d, \dots, m_d + kr$  in  $[0, N_d)$  with  $m_d + ir = n_d$  such that  $(m_d, r)$  is both perfect of order  $i$  and totally saturated. Then most elements of  $[N_d/3, 2N_d/3]$  are good, or more precisely*

$$|\{n_d \in [N_d/3, 2N_d/3] : n_d \text{ is not good}\}| \leq \frac{1}{100kF_0(L_{d-1})} N_d.$$

**Proof** This will be a simple counting argument. There will be a slight technical difficulty in that some of the error terms have size equal to some power of  $\varepsilon_{d-1}$ , which is somewhat large compared to the factors of  $\delta_{c_{d-1}(\infty)}$  which appear in (26), but by counting things carefully (and taking advantage of the upper bound (25) as well as the lower bound (26)) the factors of  $\delta_{c_{d-1}(\infty)}$  will eventually cancel each other out harmlessly.

First of all, observe from (26) that

$$\begin{aligned} |\{n_d \in [N_d/3, 2N_d/3] : |\{(m_d, r) \text{ perfect of order } i : m_d + ir = n_d\}| < (100k)^{-i-1} \delta_{c_{d-1}(\infty)}^i N_d\}| \\ < i\varepsilon_{d-1} \delta_{c_{d-1}(\infty)}^{k+1} N_d. \end{aligned}$$

Also, since  $\Omega \subset [i+1, k]$ , the component  $P[(\lambda_j)_{j \in \Omega}]_{i, n_d}$  does not depend on  $(\lambda_j)$ , so by (27)

$$|\{n_d \in [N_d/3, 2N_d/3] : P[(\lambda_j)_{j \in \Omega}]_{i, n_d} \text{ is not saturated for some } (\lambda_j)_{j \in \Omega} \in [0, L_{d-1})^\Omega\}| \leq 2\varepsilon_{d-1} N_d.$$

Both these terms are acceptable by (19). Thus it will suffice to prove that

$$|X| \leq \frac{1}{200kF_0(L_{d-1})} N_d$$

(for instance), where  $X$  is the set of those  $n_d \in [N_d/3, 2N_d/3]$  such that  $n_d$  is not good, but

$$|\{(m_d, r) \text{ perfect of order } i : m_d + ir = n_d\}| \geq (100k)^{-i-1} \delta_{c_{d-1}(\infty)}^i N_d \quad (28)$$

and such that  $P[(\lambda_j)_{j \in \Omega}]_{i, n_d}$  is saturated for all  $(\lambda_j)_{j \in \Omega} \in [0, F_0(L_{d-1}))^\Omega$ .

Now suppose that  $n_d \in X$ . From (28) we have  $(100k)^{-i-1} \delta_{c_{d-1}(\infty)}^i N_d$  pairs  $(m_d, r)$   $m_d + jr = n_d$  and  $(m_d, r)$  perfect of order  $i$ . Since  $n_d$  is not good, we know that for each such pair  $(m_d, r)$ , at least one  $\Omega$ -tuple  $(\lambda_j)_{j \in \Omega} \in [0, F_0(L_{d-1}))^\Omega$  exists such that  $P[(\lambda_j)_{j \in \Omega}]$  is not totally saturated. Since  $(m_d, r)$  is perfect of order  $i$ , we know that the first  $i+1$  components  $P[(\lambda_j)_{j \in \Omega}]_0, \dots, P[(\lambda_j)_{j \in \Omega}]_{i, m_d+ir}$  are already saturated, so we must have  $P[(\lambda_j)_{j \in \Omega}]_{i', m_d+i'r}$  unsaturated for some  $i < i' \leq k$ . Thus

$$\begin{aligned} \sum_{(m_d, r) \text{ perfect of order } i : m_d + jr = n_d} \sum_{(\lambda_j)_{j \in \Omega} \in [0, L_{d-1})^\Omega} \sum_{i < i' \leq k} \\ \mathbf{1}_{P[(\lambda_j)_{j \in \Omega}]_{i', m_d+i'r} \text{ unsaturated}} \geq (100k)^{-i-1} \delta_{c_{d-1}(\infty)}^i N_d. \end{aligned}$$



Summing over all  $n_d \in X$ , we obtain

$$\sum_{(m_d, r) \text{ perfect of order } i} \sum_{(\lambda_j)_{j \in \Omega} \in [0, L_{d-1}]^\Omega} \sum_{i < i' \leq k} 1_{P[(\lambda_j)_{j \in \Omega}]_{i', m_d + i'r} \text{ unsaturated}} \geq (100k)^{-i-1} \delta_{c_{d-1}(\infty)}^i |X| N_d.$$

The number of possible  $(\lambda_j)_{j \in \Omega}$  is at most  $F_0(L_{d-1})^{k+1}$ , and the number of  $i'$  is at most  $k$ , so by the pigeonhole principle we can find a  $(\lambda_j)_{j \in \Omega}$  and  $i < i' \leq k$  such that

$$\sum_{(m_d, r) \text{ perfect of order } i} 1_{P[(\lambda_j)_{j \in \Omega}]_{i', m_d + i'r} \text{ unsaturated}} \geq \frac{1}{k F_0(L_{d-1})^{k+1}} (100k)^{-i-1} \delta_{c_{d-1}(\infty)}^i |X| N_d.$$

We rearrange this using the substitution  $n'_d := m_d + i'r$  as

$$\begin{aligned} & \sum_{n'_d \in [0, N_d]: P[(\lambda_j)_{j \in \Omega}]_{i', n'_d} \text{ unsaturated}} |\{(m_d, r) \text{ perfect of order } i : m_d + i'r = n'_d\}| \\ & \geq \frac{1}{k F_0(L_{d-1})^{k+1}} (100k)^{-i-1} \delta_{c_{d-1}(\infty)}^i |X| N_d \end{aligned}$$

or (by definition of  $f$ )

$$\sum_{n'_d \in [0, N_d]: P[(\lambda_j)_{j \in \Omega}]_{i', n'_d} \text{ unsaturated}} f_{P[(\lambda_j)_{j \in \Omega}], i, i'}(n'_d) \geq \frac{1}{k F_0(L_{d-1})^{k+1}} (100k)^{-i-1} \delta_{c_{d-1}(\infty)}^i |X| N_d. \quad (29)$$

On the other hand, from (25) and (24) we have

$$\sum_{n'_d \in [0, N_d]: f_{P[(\lambda_j)_{j \in \Omega}], i, i'}(n'_d) > (100k)^i \delta_{c_{d-1}(\infty)}^i N_d} f_{P[(\lambda_j)_{j \in \Omega}], i, i'}(n'_d) \leq i \varepsilon_{d-1} \delta_{c_{d-1}(\infty)}^{k+1} N_d^2$$

while from (27) we have

$$\begin{aligned} & \sum_{n'_d \in [0, N_d]: f_{P[(\lambda_j)_{j \in \Omega}], i, i'}(n'_d) \leq (100k)^i \delta_{c_{d-1}(\infty)}^i N_d; P[(\lambda_j)_{j \in \Omega}]_{i', n'_d} \text{ unsaturated}} f_{P[(\lambda_j)_{j \in \Omega}], i, i'}(n'_d) \\ & \leq (100k)^i \delta_{c_{d-1}(\infty)}^i N_d \times 2\varepsilon_{d-1} N_d \end{aligned}$$

so on adding we have

$$\sum_{n'_d \in [0, N_d]: P[(\lambda_j)_{j \in \Omega}]_{i', n'_d} \text{ unsaturated}} f_{P[(\lambda_j)_{j \in \Omega}], i, i'}(n'_d) \leq 4(100k)^i \delta_{c_{d-1}(\infty)}^i \varepsilon_{d-1} N_d^2$$

(since  $\varepsilon_{d-1}$  and  $\delta_{c_{d-1}(\infty)}$  are at most 1). Combining this with (29) we obtain the bound

$$|X| \leq 4k(100k)^{2i+1} L_{d-1}^{k+1} \varepsilon_{d-1} N_d;$$

the claim then follows by (19) (and the trivial bound  $i \leq k$ ).  $\blacksquare$

Combining this lemma with Lemma 1.2, we can find an increasing arithmetic progression  $\{n_d + \lambda_i s : \lambda_i \in [0, F_0(L_d)]\}$  of length  $F_0(L_d)$  in  $[N_d/3, 2N_d/3]$ , all of whose elements are good. Thus for each  $\lambda_i \in [0, F_0(L_d)]$ , there exists a progression  $m_d(\lambda_i), \dots, m_d(\lambda_i) + kr(\lambda_i)$  in  $[0, N_d]$  such that  $m_d(\lambda_i) + ir(\lambda_i) = n_d + \lambda_i s$  which is

perfect of order  $i$  and which is totally saturated. Thus if we define the progressions  $Q[(\lambda_j)_{j \in \Omega \cup \{i\}}]$  for any  $(\lambda_j)_{j \in \Omega \cup \{i\}} \in [0, F_0(L_d)]^{\Omega \cup \{i\}}$  by the formula

$$Q[(\lambda_j)_{j \in \Omega \cup \{i\}}] = \bigcup_{i'=0}^k P[(\lambda_j)_{j \in \Omega}]_{i', m_d(\lambda_i) + i' r(\lambda_i)}$$

we see that  $Q[(\lambda_j)_{j \in \Omega \cup \{i\}}]$  is a progression of rank  $d$  and dimensions  $(N_1, \dots, N_{d-1}, k+1)$  which is completely saturated and perfect of order  $i$ . In particular, the color of the first  $i$  components of  $Q[(\lambda_j)_{j \in \Omega \cup \{i\}}]$  are the perfect color  $c_{d-1}$ , which clearly does not depend on any of the  $\lambda_j$  including  $\lambda_i$ . As for the remaining components, the independence properties of these components from the  $\lambda_i$  follow from the corresponding independence properties of the corresponding components of  $P[(\lambda_j)_{j \in \Omega}]$  (note that the quantities  $m_d(\lambda_i)$  and  $r(\lambda_i)$  do not depend on any of the  $\lambda_j$  for any  $j \in \Omega$ ). Also if we fix all the  $\lambda_j$  for  $j \in \Omega$ , then

$$Q[(\lambda_j)_{j \in \Omega \cup \{i\}}]_i = P[(\lambda_j)_{j \in \Omega}]_{i, m_d(\lambda_i) + i r(\lambda_i)} = P[(\lambda_j)_{j \in \Omega}]_{i, n_d + \lambda_i}$$

which forms an increasing arithmetic progression in  $\lambda_i$ . The analogous arithmetic structural properties for the other parameters  $\lambda_j$  again follow from the corresponding properties for  $P$ . Thus we have constructed a well-arranged sequence of progressions  $Q[]$  of rank  $d$  and parameterization  $\Omega \cup \{i\}$ , as desired.  $\blacksquare$

## 7. PROOF OF PROPOSITION 5.5

We now turn to Proposition 5.5, which is the harder of the two propositions. The main task is to prove the Proposition in the case  $\Omega = \{i\}$ , in which case we let  $Q$  simply be one of the elements of the sequence. In other words, we will show

**Proposition 7.1.** *Let  $d \geq 1$  and  $0 \leq i \leq k$ . Suppose that we have a well-arranged sequence of progressions  $P[]$  of rank  $d+1$  and parameterization  $\{i\}$  such that  $P[\lambda_i]$  is homogeneous of order  $i$  for all  $\lambda_i \in [0, F_0(L_d)]$ . Then there exists  $\lambda_i \in [0, F_0(L_d)]$  such that  $P[\lambda_i]$  is homogeneous of order  $i+1$ .*

This proposition clearly gives Proposition 5.5 in the case  $\Omega = \{i\}$ . To handle the general case, when  $\Omega$  contains elements other than  $i$ , we simply fix all the variables  $\lambda_j$  for  $j \in \Omega \setminus \{i\}$  and apply the above proposition to locate a  $\lambda_i$  such that  $P[(\lambda_j)_{j \in \Omega}]$  is homogeneous of order  $i+1$ . Note that the concept of being homogeneous of order  $i+1$  depends only on the colors of the first  $i+1$  components of the progression under consideration, which by the independence property (ii) of a well-arranged sequence, does not depend on any of the  $\lambda_j$  for  $j \in \Omega \setminus \{i\}$ . Thus we can choose  $\lambda_i$  independently of all the other  $\lambda_j$  in such a way that  $P[(\lambda_j)_{j \in \Omega}]$  is always homogeneous of order  $i+1$ . We can then define  $Q[]$  by restricting  $P$  to this value of  $\lambda_i$ , thus i.e.

$$Q[(\lambda_j)_{j \in \Omega \setminus \{i\}}] := P[(\lambda_j)_{j \in \Omega}].$$

The complete saturation, independence, and arithmetic structure properties of  $Q$  then trivially follow from those of  $P$ , and the claim follows.

It remains to prove Proposition 7.1. Let us first state what we are trying to achieve. From (25), (26) we already have

$$|\{n_d \in [0, N_d] : f_{P[\lambda_i], i, j}(n_d) > (100k)^i \delta_{c_{d-1}}(\infty)^i N_d\}| \leq i \varepsilon_{d-1} \delta_{c_{d-1}}(\infty)^{k+1} N_d \quad (30)$$

and

$$|\{n_d \in [N_d/3, 2N_d/3] : f_{P[\lambda_i], i, j}(n_d) < (100k)^{-i-1} \delta_{c_{d-1}}(\infty)^i N_d\}| \leq i \varepsilon_{d-1} \delta_{c_{d-1}}(\infty)^{k+1} N_d \quad (31)$$

for all  $\lambda_i \in [0, F_0(L_d)]$  (in fact the left-hand sides here are independent of  $\lambda_i$  thanks to the independence property of well-arranged sequences) and all  $i \leq j \leq k$ , and we are trying to find a *single*  $\lambda_i \in [0, F_0(L_d)]$  such that

$$\begin{aligned} |\{n_d \in [0, N_d] : f_{P[\lambda_i], i+1, j}(n_d) > (100k)^{i+1} \delta_{c_{d-1}}(\infty)^{i+1} N_d\}| &\leq (i+1) \varepsilon_{d-1} \delta_{c_{d-1}}(\infty)^{k+1} N_d \\ |\{n_d \in [N_d/3, 2N_d/3] : f_{P[\lambda_i], i+1, j}(n_d) < (100k)^{-i-2} \delta_{c_{d-1}}(\infty)^i N_d\}| &\leq (i+1) \varepsilon_{d-1} \delta_{c_{d-1}}(\infty)^{k+1} N_d \end{aligned}$$

for all  $i < j \leq k$  (One does not need to verify that  $P[\lambda_i]$  is completely saturated, as this is given to us by the hypothesis that  $P[\lambda_i]$  is already homogeneous of order  $i$ ). We introduce the small quantity<sup>1</sup>

$$\varepsilon := \varepsilon_{d-1}^{100k} \delta_{c_{d-1}}(\infty)^{100k} \quad (32)$$

and observe that to prove our claim it will suffice to find a single  $\lambda_i \in [0, F_0(L_d)]$  such that

$$f_{P[\lambda_i], i+1, j}(n_d) = \delta_{c_{d-1}}(\infty) f_{P[\lambda_i], i, j}(n_d) + O_k(\varepsilon N_d) \text{ for all } i < j \leq k \quad (33)$$

for all but  $O_k(\varepsilon)N_d$  many values of  $n_d \in [0, N_d]$ . To explain why one could hope for such a statement, let us rephrase the task in terms of graph theory. We temporarily fix  $i < j \leq k$ , and introduce a bipartite graph  $G_{ij}$ , connecting the interval  $A_i := [0, N_d]$  to another copy of the interval  $A_j := [0, N_d]$  (which strictly speaking we should label differently to emphasize the bipartite nature of the graph, but let us ignore this technicality), thus the edge set  $E(G_{ij})$  can be thought of as a subset of  $A_i \times A_j$ . This graph is defined as follows: we connect an element  $a_i \in A_i$  with an element  $a_j \in A_j$  by an edge if we can find a progression  $m_d, m_d + s, \dots, m_d + ks$  in  $[0, N_d]$  which is perfect of order  $i$  (i.e. all the components  $P[\lambda_i]_{i', m_d+i's}$  have the perfect color for all  $0 \leq i' < i$ ; note that the choice of  $\lambda_i$  is not relevant here by the independence property of well-arranged sequences), and such that  $m_d + is = a_i$  and  $m_d + js = a_j$ . Note that each pair  $(a_i, a_j)$  is associated to at most one progression  $m_d, \dots, m_d + ks$  and so the graph has no multiplicity. For each  $\lambda_i$ , we let  $A_i[\lambda_i] \subset A_i$  be the set of those  $a_i \in A_i$  for which  $P[\lambda_i]_{i, a_i}$  has the perfect color. Unraveling the definitions, we see that

- $f_{P[\lambda_i], i, j}(a_j)$  (which does not depend on  $\lambda_i$ ) is precisely the number of edges in  $G_{ij}$  connecting an element of  $A_i$  to  $a_j$ ; and

<sup>1</sup>At this point it is useful to note the relative magnitudes here are basically

$$1/\varepsilon_{d-1} \ll N_{d-1} \ll 1/\delta_{c_{d-1}}(\infty) \ll 1/\varepsilon \ll L_d \ll F_0(L_d) \ll 1/\varepsilon_d \ll N_d$$

thanks to (16), (19), (20), (5); actually, the quantity  $1/\delta_{c_{d-1}}(\infty)$  could be unexpectedly smaller than what (5) would suggest, but that will help us, not hurt us. As discussed earlier, the basic point is that  $L_d$  and  $F_0(L_d)$  are much smaller than  $N_d$ , but much larger than  $1/\varepsilon$  or any quantity subscripted by  $d-1$ .

- $f_{P[\lambda_i], i+1, j}(a_j)$  is precisely the number of edges in  $G_{ij}$  connecting an element of  $A_i[\lambda_i]$  to  $a_j$ .

Furthermore, since  $P[\lambda_i]$  is completely saturated, the progression  $P[\lambda_i]_i$  obeys (9) (with one smaller rank), and hence

$$|A_i[\lambda_i]| = (\overline{\delta_{c_{d-1}}}(\infty) + O(\varepsilon_d))N_d, \quad (34)$$

i.e.  $A_i[\lambda_i]$  has density roughly  $\overline{\delta_{c_{d-1}}}(\infty)$  in  $A_i$ . In particular we would also expect  $f_{P[\lambda_i], i+1, j}(a_j)$  to be roughly  $\overline{\delta_{c_{d-1}}}(\infty)$  times as large as  $f_{P[\lambda_i], i, j}(a_j)$ , which is (33).

Unfortunately we are not done yet, because the graph  $G_{ij}$  and the set  $A_i[\lambda_i]$  could “conspire” to create multiplicities  $f_{P[\lambda_i], i+1, j}(a_j)$  which are quite different from what one would expect given the multiplicities of  $G_{ij}$  and the density of  $A_i[\lambda_i]$ . Indeed, for any *individual*  $\lambda_i$ , one could easily make (33) fail. However, it will turn out (from an application of the Szemerédi regularity lemma) that in order to make (33) fail, one must make the relative density  $A_i[\lambda_i]$  on some component  $A_i^{m;ij}$  of  $A_i$  to be anomalously large or anomalously small. If one does this for each  $\lambda_i$ , then using van der Waerden’s theorem we will eventually be able to make a *single* component  $A_i^{m;ij}$  consistently anomalous (either anomalously large or anomalously small) for *all*  $\lambda_i$  in a long progression (of length  $L_d$  now rather than  $F_0(L_d)$ ). But then by some density counting arguments we will be able obtain a contradiction from these facts, (34), and the fact (from (17)) that progressions of  $L_d$  cannot have too large a color density.

We turn to the details. We begin by counting the number of edges in  $G_{ij}$ .

**Lemma 7.2.** *The number  $|E(G_{ij})|$  of edges in  $G_{ij}$  obeys the estimates*

$$(100k)^{-i-2}\delta_{c_{d-1}}(\infty)^i N_d^2 \leq |E(G_{ij})| \leq (100k)^{i+1}\delta_{c_{d-1}}(\infty)^i N_d^2.$$

*In other words, up to factors depending only on  $k$ , the edge density of  $G_{ij}$  is comparable to  $\delta_{c_{d-1}}(\infty)^i$ .*

**Proof** Observe that

$$|E(G_{ij})| = \sum_{a_j \in [0, N_d)} f_{P[\lambda_i], i, j}(a_j).$$

Applying (30) and the crude bound  $f_{P[\lambda_i], i, j}(a_j) \leq N$  to handle the exceptional set, we obtain the desired upper bound. For the lower bound, we trivially bound

$$|E(G_{ij})| \geq (100k)^{-i-1}\delta_{c_{d-1}}(\infty)^i N_d |\{a_j \in [N_d/3, 2N_d/3) : f_{P[\lambda_i], i, j}(a_j) \geq (100k)^{-i-1}\delta_{c_{d-1}}(\infty)^i N_d\}|$$

and use (31). ■

Now we use the Szemerédi regularity lemma. We shall use the modern version of this lemma (see e.g. my expository note [5], or many other references in the literature), rather than the original one in [3]. For any non-empty subsets  $A'_i \subseteq A_i$ ,  $A'_j \subseteq A_j$ , define the density  $0 \leq d(A'_i, A'_j) \leq 1$  of this pair to be the quantity

$$d(A'_i, A'_j) := \frac{|E(G_{ij}) \cap (A'_i \times A'_j)|}{|A'_i \times A'_j|}.$$

We say that a pair  $A'_i, A'_j$  is  $\epsilon^2$ -regular if we have

$$|E(G_{ij}) \cap (A''_i \times A''_j)| = d(A'_i, A'_j)|A''_i \times A''_j| + O(\epsilon^2)|A'_i \times A'_j|$$

for all subsets  $A''_i \subseteq A'_i, A''_j \subseteq A'_j$ . The Szemerédi regularity lemma asserts that there is a constant  $C(\epsilon)$  depending only on  $\epsilon$ , and partitions<sup>2</sup>

$$A_l = A_l^{1;ij} \cup \dots \cup A_l^{M_{ij};ij} \cup A_l^{error;ij} \text{ for } l = i, j$$

with  $M_{ij} \leq C(\epsilon)$  and  $|A_l^{M_{ij};ij}| = O(\epsilon)N_d$  and  $|A_l^{m;ij}| \sim N_d/M_{ij}$  for  $m \in [1, M_{ij})$ , such that all but  $O(\epsilon^4 M_{ij}^2)$  pairs  $(A_i^{m;ij}, A_j^{m';ij})$  are  $\epsilon^2$ -regular. Actually, we can refine the last statement a bit, to say that for every  $m'$ , the pair  $(A_i^{m;ij}, A_j^{m';ij})$  is epsilon regular for all but  $O(\epsilon^2 M_{ij})$  values of  $m$ , since the values of  $m'$  which do not obey such a property can be safely absorbed into the exceptional set (and then one has to reduce  $M_{ij}$  accordingly, e.g. by removing the corresponding values of  $m$  into an exceptional set also).

Let us fix this decomposition. The point of this decomposition is that if  $a_i \in A_i^{m;ij}$  and  $a_j \in A_j^{m';ij}$ , then one should behave as if  $a_i$  were connected to  $a_j$  in  $G$  with “probability”  $d(A_i^{m;ij}, A_j^{m';ij})$ . In particular, one now expects to have the approximation

$$f_{P[\lambda_i, i+1, j]}(a_j) \approx \bar{f}_{P[\lambda_i, i+1, j]}(a_j)$$

where

$$\bar{f}_{P[\lambda_i, i+1, j]}(a_j) := \sum_{m, m' \in [1, M_{ij}]} d(A_i^{m;ij}, A_j^{m';ij}) |A_i[\lambda_i] \cap A_i^{m;ij}| \mathbf{1}_{a_j \in A_j^{m';ij}}. \quad (35)$$

Similarly one expects

$$f_{P[\lambda_i, i, j]}(a_j) \approx \bar{f}_{P[\lambda_i, i, j]}(a_j)$$

where

$$\bar{f}_{P[\lambda_i, i, j]}(a_j) := \sum_{m, m' \in [1, M_{ij}]} d(A_i^{m;ij}, A_j^{m';ij}) |A_i^{m;ij}| \mathbf{1}_{a_j \in A_j^{m';ij}}. \quad (36)$$

The point is that the dependence of  $\bar{f}_{P[\lambda_i, i+1, j]}$  on  $\lambda_i$  is much more controllable than the original function  $f_{P[\lambda_i, i+1, j]}(a_j)$ , as one only needs to control the  $M_{ij}$  numbers  $|A_i[\lambda_i] \cap A_i^{m;ij}|$ . The functions  $f_{P[\lambda_i, i, j]}(a_j), \bar{f}_{P[\lambda_i, i, j]}(a_j)$ , of course, do not depend on  $\lambda_i$  at all.

Let us first verify that  $f_{P[\lambda_i, i+1, j]}$  and  $\bar{f}_{P[\lambda_i, i+1, j]}$  are in fact quite close (and similarly with  $i+1$  replaced by  $i$ ).

**Proposition 7.3.** *We have*

$$f_{P[\lambda_i, i+1, j]}(a_j) = \bar{f}_{P[\lambda_i, i+1, j]}(a_j) + O(\epsilon N_d) \quad (37)$$

for all  $a_j \in [0, N_d]$  with at most  $O(\epsilon N_d)$  exceptions. Similarly we have

$$f_{P[\lambda_i, i, j]}(a_j) = \bar{f}_{P[\lambda_i, i, j]}(a_j) + O(\epsilon N_d) \quad (38)$$

<sup>2</sup>The notation here is unfortunately a bit heavy, but we have to keep track of the dependence on  $j$  as the definition of homogeneity will eventually require us to let  $j$  range over all values in  $(i, k]$ . However for the next few paragraphs one can suppress the dependence on  $j$ .

for all  $a_j \in [0, N_d)$  with at most  $O(\epsilon N_d)$  exceptions.

**Proof** We just prove (37), as the proof of (38) is the same. The entire case  $a_j \in A_j^{error;ij}$  can be placed in the exceptional set, so now let us assume that  $a_j \in A_j^{m';ij}$  for some  $m'$ ; for this fixed value of  $m'$  we will show that (37) holds for all  $a_j \in A_j^{m';ij}$  with at most  $O(\epsilon N_d/M_{ij})$  exceptions.

For  $a_j \in A_j^{m';ij}$ , we have

$$f_{P[\lambda_i], i+1, j}(a_j) = \sum_{m \in [1, M_{ij}]} |\{a_i \in A_i^{m;ij} : (a_i, a_j) \in E(G)\}| + |\{a_i \in A_i^{error;ij} : (a_i, a_j) \in E(G)\}|$$

and

$$\overline{f}_{P[\lambda_i], i+1, j}(a_j) = \sum_{m \in [1, M_{ij}]} d(A_i^{m;ij}, A_j^{m';ij}) |A_i[\lambda_i] \cap A_i^{m;ij}|.$$

The error term  $|\{a_i \in A_i^{error;ij} : (a_i, a_j) \in E(G)\}|$  is  $O(|A_i^{error;ij}|) = O(\epsilon N_d)$ , so it will suffice to show that

$$\sum_{m \in [1, M_{ij}]} (|\{a_i \in A_i^{m;ij} : (a_i, a_j) \in E(G)\}| - d(A_i^{m;ij}, A_j^{m';ij}) |A_i[\lambda_i] \cap A_i^{m;ij}|) = O(\epsilon N_d)$$

with at most  $O(\epsilon N_d/M_{ij})$  exceptions. By Chebyshev's inequality, it will suffice to show that

$$\sum_{a_j \in A_j^{m';ij}} \sum_{m \in [1, M_{ij}]} (|\{a_i \in A_i^{m;ij} : (a_i, a_j) \in E(G)\}| - d(A_i^{m;ij}, A_j^{m';ij}) |A_i[\lambda_i] \cap A_i^{m;ij}|) = O(\epsilon^2 N_d^2/M_{ij}).$$

The contribution of those  $m$  for which  $(A_i^{m;ij}, A_j^{m';ij})$  is not  $\epsilon^2$ -regular is acceptable since each such  $m$  contributes at most  $O(N_d^2/M_{ij}^2)$  and there are only  $O(\epsilon^2 M_{ij})$  such  $m$ . So we can restrict to the  $\epsilon^2$ -regular pairs. By the triangle inequality, it then will suffice to show that

$$\sum_{a_j \in A_j^{m';ij}} \sum_{m \in [1, M_{ij}]} (|\{a_i \in A_i^{m;ij} : (a_i, a_j) \in E(G)\}| - d(A_i^{m;ij}, A_j^{m';ij}) |A_i[\lambda_i] \cap A_i^{m;ij}|) = O(\epsilon^2 N_d^2/M_{ij}^2)$$

whenever  $(A_i^{m;ij}, A_j^{m';ij})$  is  $\epsilon^2$ -regular. Without the absolute values, the summand has mean zero, so it suffices to show that

$$\sum_{a_j \in X} |\{a_i \in A_i^{m;ij} : (a_i, a_j) \in E(G)\}| - d(A_i^{m;ij}, A_j^{m';ij}) |A_i[\lambda_i] \cap A_i^{m;ij}| = O(\epsilon^2 N_d^2/M_{ij}^2)$$

uniformly for all subsets  $X$  of  $A_j^{m';ij}$ . But the left-hand side is the same as

$$|E(G) \cap ((A_i[\lambda_i] \cap A_i^{m;ij}) \times X)| - d(A_i^{m;ij}, A_j^{m';ij}) |A_i[\lambda_i] \cap A_i^{m;ij}| |X| = O(\epsilon^2 N_d^2/M_{ij}^2)$$

and the claim now follows by  $\epsilon^2$ -regularity.  $\blacksquare$

In light of (37), (38), we see that in order to prove (33) it will suffice to find a  $\lambda_i \in [0, F_0(L_d))$  such that

$$\overline{f_{P[\lambda_i], i+1, j}(a_j)}} = \delta_{c_{d-1}}(\infty) \overline{f_{P[\lambda_i], i, j}(a_j)}} + O_k(\epsilon N_d) \text{ for all } i < j \leq k, a_j \in [0, N_d).$$

On the other hand, from (35), (36), and the triangle inequality we have

$$\begin{aligned} |\overline{f_{P[\lambda_i, i+1, j]}(a_j)} - \delta_{c_{d-1}}(\infty) \overline{f_{P[\lambda_i, i, j]}(a_j)}| &\leq \sup_{m' \in [1, M_{ij}]} \sum_{m \in [1, M_{ij}]} d(A_i^{m; ij}, A_j^{m'; ij}) \\ &\quad \|\overline{A_i[\lambda_i] \cap A_i^{m; ij}} - \delta_{c_{d-1}}(\infty) \overline{A_i^{m; ij}}\| \\ &\leq M_{ij} \sup_{m, m' \in [1, M_{ij}]} \|\overline{A_i[\lambda_i] \cap A_i^{m; ij}} - \delta_{c_{d-1}}(\infty) \overline{A_i^{m; ij}}\| \end{aligned}$$

and so it suffices to find a  $\lambda_i \in [0, F_0(L_d))$  such that

$$|A_i[\lambda_i] \cap A_i^{m; ij}| = \delta_{c_{d-1}}(\infty) |A_i^{m; ij}| + O_k(\epsilon N_d / M_{ij}) \text{ for all } i < j \leq k; m \in [1, M_{ij}].$$

The point here is that the free parameter  $a_j$ , which ranged over an enormous set (of size  $N_d$ , which is larger than all the other quantities currently in play, and in particular is significantly larger than the number of possible  $\lambda_i$  at our disposal, which is  $F_0(L_d)$ ) has now been replaced by the parameters  $m, m'$ , which range over a much smaller set (of size  $M_{ij}$ , which as it depends ultimately on  $\epsilon_{d-1}$  and  $\delta_{c_{d-1}}(\infty)$  will in particular be quite small compared to  $F_0(L_d)$ ). This makes our task far easier to perform (we now are trying to locate a solution to a problem which has many more degrees of freedom than constraints, whereas previously the situation was reversed). This is the power of the Szemerédi regularity lemma: to reduce a complicated object (a graph on many vertices) to a simpler one (a probabilistic graph on a much smaller set of vertices).

We still have to locate  $\lambda_i$ . Suppose for contradiction that this was not possible. Then for each  $\lambda_i$  in  $[0, F_0(L_d))$ , there exists  $1 < j \leq k$ ,  $m \in [1, M_{ij}]$ , and a sign  $\pm$  such that  $A_i[\lambda_i]$  has anomalous density in  $A_i^{m; ij}$ :

$$\pm (|A_i[\lambda_i] \cap A_i^{m; ij}| - \delta_{c_{d-1}}(\infty) |A_i^{m; ij}|) \geq \epsilon N_d / M_{ij}. \quad (39)$$

At present the quantities  $j, m, \pm$  depend on  $\lambda_i$ . However, we now invoke our second powerful tool, the van der Waerden theorem, to eliminate this dependence (while still retaining arithmetic structure on the  $\lambda_i$ ). From the Szemerédi regularity lemma we had  $M_{ij} \leq C(\epsilon)$ . Combining this with (32), (16), (20) we have the number of triplets  $(j, m, \pm)$  is bounded by

$$\leq 2kC(\epsilon) = C(\epsilon_{d-1}, \delta_{c_{d-1}}(\infty)) = C(N_{d-1}) \ll L_d$$

(if  $F_0$  is sufficiently fast) and so the map  $\lambda_i \rightarrow (j, m, \pm)$  colors the interval  $[0, F_0(L_d))$  into at most  $L_d$  colors. By the van der Waerden theorem we can thus find (if  $F_0$  is sufficiently fast growing<sup>3</sup>) an increasing arithmetic progression  $R$  in  $[0, F_0(L_d))$  of length  $L_d$  and *fixed*  $j, m, \pm$  such that (39) holds for all  $\lambda_i \in R$ ; thus we have a consistently anomalous density in this progression  $R$ . We now use some density counting arguments to obtain a contradiction from this and (17).

---

<sup>3</sup>This step is the step that causes  $F_0$  to grow extremely fast; even the Szemerédi regularity lemma only demands tower-type growth of  $F_0$ , whereas the best elementary bound on van der Waerden's theorem, due to Shelah, is a triply iterated tower. This in turn leads to very poor final bounds on Szemerédi's theorem. The bound of Gowers gives much better bounds here - double exponential, to be precise - but it would be rather perverse to use Gowers' bound here, since the argument in [2] gives Szemerédi's theorem directly.

Summing (39) for all  $\lambda_i \in R$  (and taking full advantage of the consistency of the parameters  $j, m, \pm$ ) we see that

$$\left| \sum_{\lambda_i \in R} |A_i[\lambda_i] \cap A_i^{m;ij} - \delta_{c_{d-1}}(\infty)| |A_i^{m;ij}| L_d \right| \geq \epsilon L_d N_d / M_{ij}. \quad (40)$$

On the other hand, by applying (34) for all  $\lambda_i \in R$  and summing to obtain

$$\sum_{\lambda_i \in R} |A_i[\lambda_i]| = (\overline{\delta_{c_{d-1}}}(\infty) + O(\epsilon_d)) L_d N_d. \quad (41)$$

To obtain a contradiction we need to somehow restrict (41) to  $A_i^{m;ij}$ . As we shall see, this will be possible thanks to the control (17) on the perfect color density at scale  $L_d$ .

We turn to the details. We rearrange the left-hand side to obtain

$$\sum_{a_i \in [0, N_d]} |\{\lambda_i \in R : P[\lambda_i]_{i, a_i} \text{ has the perfect color}\}| = (\overline{\delta_{c_{d-1}}}(\infty) + O(\epsilon_d)) L_d N_d.$$

Note that the progression  $Q_{a_i} := \bigcup_{\lambda_i \in R} P[\lambda_i]_{i, a_i}$  is an increasing progression of rank  $d$  and dimensions  $(N_1, \dots, N_{d-1}, L_d)$ . Thus we can rearrange the left-hand side again to obtain

$$\sum_{a_i \in [0, N_d]} \delta_{c_{d-1}}(Q_{a_i}) = (\overline{\delta_{c_{d-1}}}(\infty) + O(\epsilon_d)) N_d.$$

We now apply

**Lemma 7.4.** *We have*

$$\sum_{a_i \in [0, N_d]} (\delta_{c_{d-1}}(Q_{a_i}) - \overline{\delta_{c_{d-1}}}(\infty))_+ = O\left(\frac{1}{F_0(N_1 \dots N_{d-1})}\right) N_d$$

**Proof** There are two contributions, depending on whether  $Q_{a_i}$  obeys the condition

$$|\sigma_{N^{(d-1)}}(Q_{a_i}) - \overline{\sigma_{N^{(d-1)}}}(\infty)| \leq \sqrt{\mu_{d-1}(L_d)}$$

or not. If it does, then by (3) we have

$$\delta_{c_{d-1}}(Q_{a_i}) \leq \overline{\delta_{c_{d-1}}}(L_d)$$

and this contribution will be acceptable from (17). In the other case, we bound the summand crudely by 1, and reduce to showing that

$$|\{a_i \in [0, N_d] : |\sigma_{N^{(d-1)}}(Q_{a_i}) - \overline{\sigma_{N^{(d-1)}}}(\infty)| > \sqrt{\mu_{d-1}(L_d)}\}| = O\left(\frac{1}{F_0(N_1 \dots N_{d-1})}\right) N_d.$$

To see this, we first observe from (8), (7) and the saturation of the  $P[\lambda_i]_i$  that

$$|\{a_i \in [0, N_d] : P[\lambda_i]_{i, a_i} \text{ is saturated}\}| = (\overline{\sigma_{N^{(d-1)}}}(\infty) + O(\epsilon_d^4)) N_d$$

for all  $\lambda_i \in R$ ; averaging this in  $\lambda_i$  and rearranging the left-hand side we obtain

$$\sum_{a_i \in [0, N_d]} \sigma_{N^{(d-1)}}(Q_{a_i}) = (\overline{\sigma_{N^{(d-1)}}}(\infty) + O(\epsilon_d^4)) N_d$$

and thus

$$\sum_{a_i \in [0, N_d]} (\sigma_{N^{(d-1)}}(Q_{a_i}) - \overline{\sigma_{N^{(d-1)}}}(\infty)) = O(\epsilon_d^4) N_d.$$



But from definition of  $\mu_{d-1}(L_d)$  we have

$$\sum_{a_i \in [0, N_d)} (\sigma_{N^{(d-1)}}(Q_{a_i}) - \overline{\sigma_{N^{(d-1)}}}(\infty))_+ = O(\mu_{d-1}(L_d))N_d$$

so by (21)

$$\sum_{a_i \in [0, N_d)} |\sigma_{N^{(d-1)}}(Q_{a_i}) - \overline{\sigma_{N^{(d-1)}}}(\infty)| = O(\mu_{d-1}(L_d))N_d$$

and the claim now follows from Chebyshev and (18). ■

From this Lemma and the preceding estimate we have (from (19), (16))

$$\sum_{a_i \in [0, N_d)} |\delta_{c_{d-1}}(Q_{a_i}) - \overline{\delta_{c_{d-1}}}(\infty)| = O\left(\frac{1}{F_0(N_1 \dots N_{d-1})}\right)N_d.$$

In particular we have

$$\sum_{a_i \in A_i^{m;ij}} |\delta_{c_{d-1}}(Q_{a_i}) - \overline{\delta_{c_{d-1}}}(\infty)| = O\left(\frac{1}{F_0(N_1 \dots N_{d-1})}\right)N_d$$

and hence

$$\sum_{a_i \in A_i^{m;ij}} \delta_{c_{d-1}}(Q_{a_i}) = \overline{\delta_{c_{d-1}}}(\infty)|A_i^{m;ij}| + O\left(\frac{1}{F_0(N_1 \dots N_{d-1})}\right)N_d.$$

We now reverse our previous manipulations, to obtain

$$\sum_{a_i \in A_i^{m;ij}} |\{\lambda_i \in R : P[\lambda_i]_{i, a_i} \text{ has the perfect color}\}| = \overline{\delta_{c_{d-1}}}(\infty)|A_i^{m;ij}| + O\left(\frac{1}{F_0(N_1 \dots N_{d-1})}\right)L_d N_d$$

which rearranges as before to 
$$\sum_{\lambda_i \in R} |A_i[\lambda_i] \cap A_i^{m;ij}| = \overline{\delta_{c_{d-1}}}(\infty)|A_i^{m;ij}| + O\left(\frac{1}{F_0(N_1 \dots N_{d-1})}\right)L_d N_d.$$

Comparing this with (40) we obtain

$$\epsilon/M_{ij} = O\left(\frac{1}{F_0(N_1 \dots N_{d-1})}\right)$$

which is a contradiction if  $F_0$  grows quickly enough, because  $\epsilon/M_{ij}$  can be bounded below by a quantity depending on  $\epsilon$ , which in turn depends on  $\epsilon_{d-1}$  and  $\overline{\delta_{c_{d-1}}}(\infty)$ , which in turn is bounded by  $N_1 \dots N_{d-1}$  thanks to (5) (and (20)). This finally completes the proof of Proposition 5.5 and hence of Szemerédi's theorem.

## REFERENCES

- [1] H. Furstenberg, *Recurrence in Ergodic theory and Combinatorial Number Theory*, Princeton University Press, Princeton NJ 1981.
- [2] T. Gowers, *A new proof of Szemerédi's theorem*, GAFA **11** (2001), 465-588.
- [3] E. Szemerédi, *On sets of integers containing no  $k$  elements in arithmetic progression*, Acta Arith. **27** (1975), 299-345.
- [4] B.L. Van der Waerden, *Beweis einer Baudetschen Vermutung*, Nieuw. Arch. Wisk. **15** (1927), 212-216.
- [5] T. Tao, *An information theoretic proof of Szemerédi's regularity lemma*, unpublished.

DEPARTMENT OF MATHEMATICS, UCLA, LOS ANGELES CA 90095-1555

*E-mail address:* `tao@math.ucla.edu`