

# Site-specific recombinatorics: *in situ* cellular barcoding with the Cre Lox system

Tom S. Weber<sup>1</sup>, Mark Dukes<sup>2</sup>, Denise Miles<sup>3</sup>, Stefan Glaser<sup>3</sup>, Shalin Naik<sup>3</sup> and Ken R. Duffy<sup>1</sup>

<sup>1</sup>Hamilton Institute, Maynooth University, Ireland

<sup>2</sup>University of Strathclyde, Glasgow, United Kingdom

<sup>3</sup>The Walter and Eliza Hall Institute of Medical Research & The University of Melbourne, Parkville, Australia

March 1, 2016

## Abstract

Cellular barcoding is a significant, recently developed, biotechnology tool that enables the familial identification of progeny of individual cells *in vivo*. Most existing approaches rely on *ex vivo* viral transduction of cells with barcodes, followed by adoptive transfer into an animal, which works well for some systems, but precludes barcoding cells in their native environment, such as those inside solid tissues. With a view to overcoming this limitation, we propose a new design for a genetic barcoding construct based on the Cre Lox system that induces randomly created stable barcodes in cells *in situ* by exploiting inherent sequence distance constraints during site-specific recombination. Leveraging this previously unused feature, we identify the cassette with maximal code diversity. This proves to be orders of magnitude higher than what is attainable with previously considered Cre Lox barcoding approaches and is well suited for its intended applications as it exceeds the number of lymphocytes or hematopoietic progenitor cells in mice. Moreover, it can be built using established technology.

**Keywords:** cell fate tracking; cellular barcoding; cre lox system; DNA stochastic programme; combinatorial explosion.

## Introduction

The fate of the progeny of two seemingly identical cells can be markedly distinct. Well studied examples include the immune system and hematopoietic system, for which the extent of clonal expansion and differentiation has been shown to vary greatly between cells of the same phenotype [1–5]. Fate and expression heterogeneity at the single-cell level are also apparent in other systems including the brain [6–8] and cancers [9–11]. Whether this heterogeneity is due to the stochastic nature of cellular decision making, reflects limitations in phenotyping, is caused by external events, or a mixture of effects, is a subject of active study [12,13]. As addressing this pivotal question through population-level analysis is not possible, tools have been developed that facilitate monitoring single cells and their offspring across generations.

Long-term fluorescence microscopy represents the most direct approach to assess fate heterogeneity at the single-cell level. Studies employing that technique are numerous [14–19], and have revealed many significant features. Filming and tracking of cell families *in vitro* remains technically challenging, is labor intensive, and only partially automatable [20,21]. Despite significant advances in the field, continuous tracking *in vivo* is confined to certain tissues, and time windows of up to twelve hours for slowly migrating cells.

A radically different approach to long-term clonal monitoring is to mark single cells with unique DNA tags via retroviral transduction, a technique known as cellular barcoding [2,11,22,23]. As tags are heritable, clonally related cells can be identified via DNA sequencing. By tagging multi-potent cells of the hematopoietic system and adoptively transferring them into irradiated mice, the contribution of single stem cells to hematopoiesis has been quantified [22,23]. Amongst other discoveries, this has revealed heterogeneity in the collection of distinct cell types produced from apparently equi-potent progenitors [4,5,24]. Current barcoding techniques are unsuitable for tagging cells *in vivo*, and typically require *ex*

*vivo* barcoding followed by adoptive cell transfer [23]. This restricts its scope to cell types such as naive lymphocytes and cancer cells, as well as hematopoietic stem and progenitors which require perturbation of the new host, usually irradiation, to enable them to engraft.

Ideally, a cellular barcoding system would inducibly mark cells in their native environment, would be non-toxic, permanent and heritable, barcodes would be easy to read with a high-throughput technique, and the system would enable labeling large numbers of cells with unique barcodes. Two recently published studies address some of these points. Sun et al. [25] employed a Dox inducible form of the Sleeping Beauty transposase to genetically tag stem cells *in situ*, and followed clonal dynamics during native hematopoiesis in mice. These tags are the random insertion site of an artificial transposon, which upon withdrawal of Dox is relatively stable. A second *in situ* cellular barcoding system based on site-specific DNA recombination with the Rci invertase has also been implemented [26–28]. Inspired by the Brainbow mouse [29], this system induces a random barcode by stochastically shuffling a synthetic cassette pre-integrated into the genome of a cell. The authors predicted high code diversity from relatively small constructs (approx. 2 kb) and demonstrated feasibility of random barcode generation in *Escherichia coli* [28].

Each of those approaches elegantly overcome shortcomings of previous systems by generating largely unique tags without significant perturbation to the system of interest, but some difficulties remain. For barcode readout, the Sleeping Beauty system requires whole-genome amplification technology and three-arm-ligation-mediated PCR to efficiently amplify unknown insertion sites. Furthermore, the random location of the transposon may impact behavior of some barcoded clones and lead to biased data. Moreover, some background transposon mobilization was detected, subverting the stability of the barcodes. The Rci invertase based system remains to be implemented outside bacteria. As with the Sleeping Beauty transposase, the method requires tight temporal control over Rci expression to make codes permanent.

Here we consider the Cre Lox system as a driver to induce *in situ* large numbers of distinct, permanent, randomly determined barcodes from a series of tightly spaced Lox sites. In contrast to the Brainbow construct [29], which relies on overlapping pairs of incompatible Lox sites recombining randomly to one of several stable DNA sequence configurations, our design exploits constraints on the distance between Lox sites that arise during DNA loop formation, a prerequisite for site-specific recombination [30–32]. This known feature has not previously been exploited, but is a crucial design element for obtaining high barcode diversity. Employing repeated usage of the same Lox site, code diversity is solely restricted by cassette size and not, as in the Brainbow construct, by the relatively small set of non-interacting Lox sites [33]. For a design without distance constraints, the maximal diversity of stable barcodes creatable with the Cre Lox system is of order  $n$ , where  $n$  is the number of Lox sites [28], but with distance constraints we establish that optimal barcode diversities of order  $n^3$  are possible. Boosting this scaling with the four incompatible Lox sites that have been reported in the literature [33] enables  $10^{12}$  distinct codes of about 600 bp each from a genetic construct as small as 2.5 kb. In combination with the CreEr system [34], this is sufficient to inducibly barcode label all naive CD8 T cells in a mouse [35] or all nucleated cells in the bone marrow [36]. Desirable features are inherently part of the Lox barcode cassette design, including: short and stable barcodes; a single barcode per cell; and robust read-out.

## Cre Lox biology

Before introducing the Lox barcode cassette, we revisit Cre Lox biology [37,38]. Cre is a bacteriophage P1 recombinase that catalyzes site-specific recombination between Lox sites. A Lox site is a 34 bp sequence composed of two 13 bp palindromic flanking regions and an asymmetric 8 bp core region (Fig. 1 A). For recombination to occur, four Cre proteins bind to the four palindromic regions of two Lox sites and form a synaptic complex. A first pair of strand exchanges leads to a Holliday junction intermediate [39]. Isomerization of the intermediate then allows a second pair of strand exchanges, and formation of the final recombinant product [32]. The DNA cleavage site is situated in the asymmetric core region. If the Lox sites are on the same chromosome, their interaction requires formation of a DNA loop. If they have the same orientation (direct repeats), recombination results in excision of the intervening sequence. If Lox sites are in the opposite orientation to each other (inverted repeats), the sequence between the sites is inverted, becoming its reverse complement (Fig. 1 B). Due to compatibility with eukaryotes, the Cre Lox system has become an essential tool in genetic engineering and a large array of transgenic mouse models with inducible cell-type specific expression of Cre have been created [34].

In *in vitro* trials with Cre mediated Lox reactions, a sharp decrease in recombination efficiency has been observed when the sequence separating two Lox sites is less than 94 bp [30]. Recombination is still detectable at low levels at 82 bp, but not at 80 bp where DNA stiffness appears to prevent DNA

loop formation, and as a consequence Lox site interaction. For the distinct, but similar, Flp/FR system this minimal distance was established to be smaller *in vivo*, with interactions possible at 74bp [31]. The existence of a minimal distance is one of the key features that we exploit to make random barcodes stable, but in our proposed design it will only prove necessary for it to be greater than 44 bp.

## Lox barcode cassettes

In complete generality, a Lox barcode cassette is a series of Lox sites interlaced with  $n$  distinguishable DNA code elements of size  $m$  bp each. On Cre expression, code elements change orientation and position, or are excised [27]. Through Cre mediated excision, the number of elements eventually decreases until reaching a stable number (Fig. 1 C). Sequences that have attained a stable number of code elements form size-stable barcodes. A cassette’s code diversity is the number of size-stable barcodes that can be generated from the cassette via site-specific recombination.

Our main result is a robust Lox cassette design that provably maximizes code diversity. The design is robust to both sequencing errors and to the minimal interaction distance between Lox sites. The analysis that leads us to the design is provided in the Optimal Design section. The identification of code element sequences that avoid misclassification due to sequencing read errors then follows. Finally, probabilistic aspects of code generation from an optimal barcode cassette are explored via Monte Carlo simulation. Lox cassettes with code elements of size 4 bp, higher order Lox interactions, and the impact of transient Cre activation, are considered in the discussion.

### A robust cassette design that maximizes code diversity

The optimal design will prove to have the orientation of both the outmost, and any two consecutive, Lox sites inverted (Fig. 1 C). Code elements between Lox sites are of size longer than four bp, but shorter than 24 bp. The lower limit ensures that elements can be chosen sufficiently distinctly to correct two sequencing errors per element. Due to the minimal Lox interaction distance, the upper limit ensures that barcodes with three code elements are size-stable.

The barcode diversity for this cassette design with  $n$  code elements under constitutive Cre expression will, as established in the Optimal Design section, transpire to be

$$\frac{(n+1)(n-1)^2}{2} + (n+1) = O(n^3), \quad (1)$$

which is maximal for code elements that are larger than four base pairs.

A good compromise between cassette, robustness to sequencing errors and barcode diversity is given by an alternating Lox cassette with 13 elements of length 7 bp each as shown in Fig. 1 C. The cassette is initially 567 bp long and generates a code diversity of 1022 barcodes. After excisions and inversions, size-stable barcodes are composed of either a single element or three elements, with lengths 75 bp and 157 bp respectively, including remaining non-interacting Lox sites. Concatenating four such cassettes with poorly-interacting Lox variants (e.g. LoxP, Lox2272, Lox5171 and m2 [33], Fig. 1 D) yields a 2268 bp construct with a size-stable code diversity of  $1022^4 \approx 10^{12}$ .

### A practical implementation

To implement Cre Lox barcoding in the mouse, one could cross mice generated from embryonic stem cells that have been transduced with the concatenated Lox barcoding cassettes described above onto a tamoxifen inducible cell-type specific CreEr expressing background [34]. An experiment is initiated by administering tamoxifen to the animal, which activates Cre and induces generation of a barcode ( $\leq 628$  bp) in each cell where Cre becomes active. Some time after activation, cells of interest are harvested and sorted for specific phenotypes, and sequenced using a next generation sequencing platform that produces read-lengths  $> 600$  bp. Cells originating from the same progenitor carry the same barcode and this information can then be used for scientific inference. To identify the frequent barcodes that are to be discarded in the analysis (see the Barcode Distribution is Heterogeneous section), in a control experiment large numbers of cells would be harvested shortly after tamoxifen administration and sequenced.

## Optimal design

A simple upper bound on the barcode diversity of  $k$  elements from a cassette initially containing  $n$  elements is the number of possible outcomes when choosing  $k$  from  $n$  elements in arbitrary order and

orientation:

$$\binom{n}{k} k! 2^k = \frac{2^k n!}{(n-k)!}.$$

Although loose, it will become clear that it captures the dominant growth,  $O(n^k)$ , indicating the importance of  $k$  in generating barcode diversity and motivating a closer look at how cassette designs influence it.

For what follows, we introduce some terminology: a cassette is alternating if the orientation of any two consecutive Lox sites is inverted (Fig. 1 C); outermost Lox sites are termed flanking Lox sites; and flanking sites are direct or inverted if they have the same or opposite orientation, respectively.

### Code diversity is determined by code element length and orientation of flanking sites

Cre recombination requires a minimal distance between the interacting Lox sites. In what follows we assume that the minimal distance for Lox interaction is 82 bp, but our results will be robust for any minimal interaction distance greater than 44 bp.

To understand how a minimal Lox-Lox interaction distance and cassette design determine size-stable barcodes and code diversity, we start with the simplest case, a barcode with a single code element (Fig. 2 A). If the code element is less than 82 bp, the barcode is size-stable irrespective of the orientation of its flanking sites. If the element is larger than 82 bp, the code is only size-stable if the flanking sites are inverted as excision will remove the element.

For a barcode with two elements, the sequence between the flanking sites contains an additional element and a Lox site (34 bp), giving a sequence of  $2m + 34$  bp. If the flanking sites have the same orientation, the barcode is size-stable if  $2m + 34 < 82$  bp, hence if  $m < 24$  bp. If they are in opposite orientation, excisions can only occur if flanking sites interact with the middle Lox site, and  $m < 82$  bp is sufficient for stability (Fig. 2 B). For given  $m$ , in general if there exists a barcode of size  $k$  with direct flanking sites, a barcode with  $k + 1$  elements is possible that has inverted flanking sites. Thus  $m$  and the orientation of the flanking sites are critical features that determine the maximum  $k$ .

In Fig. 2 C, the stability of barcodes with  $k \in \{2, 3, 4, 5\}$  is shown as a function of  $m$  for a cassette with inverted flanking sites. The stability depends on a critical distance, i.e., the largest distance between two Lox sites in the barcode that is, or can be brought into, the same orientation via recombination. As shown, barcodes of size three and four become unstable if  $m \geq 24$  bp and  $m \geq 5$  bp, respectively, while barcodes of size five or greater are always unstable.

Orientation of a cassette's flanking sites is immutable under recombination. Therefore cassettes with direct and inverted flanking sites generate barcodes with direct and inverted flanking sites only. Having seen that maximal code diversity grows as  $O(n^k)$ , and that having inverted flanking sites relative to direct ones increases the maximum size of barcodes by one, it follows that the diversity for cassettes with inverted flanking sites is of the order  $O(n^{k+1})$ . Inverted flanking sites are thus superior in terms of code diversity and are an essential design decision.

Optimality regarding the size of the elements,  $m$ , is more intricate. For  $m < 5$ , the maximum size of barcodes is four elements, and according to the formula above, their diversity grows as  $O(n^4)$ . The stability of barcodes with four elements is, however, sensitive to the minimal distance estimate (the gray interval in Fig. 2 C). In addition, the short length of code elements limits error correction, a point revisited later. Thus we focus on cassettes in the regime  $5 \text{ bp} \leq m < 24 \text{ bp}$ , which generate error-robust barcodes of up to size three and a code diversity that is insensitive to the reported minimal Lox interaction distance.

### Alternating Lox cassettes with inverted flanking sites maximize code diversity

For the orientation of the remaining Lox sites we prove, via a two-step strategy, that the alternating design produces maximal code diversity. First we derive a refined upper bound for the diversity that takes into account the structure of the Lox cassette, but ignores constraints imposed by the recombination process. We then show that alternating Lox cassettes with inverted flanking sites and  $n \geq 7$  elements are unconstrained in terms of barcode generation via sequential recombination events, thus achieving this upper bound.

## An upper bound for Lox barcode diversity

During Cre induced recombination, Cre proteins cleave the core region of the interacting Lox sites asymmetrically [32]. The sequences between subsequent cleavage sites are not affected by Cre and represent the fundamental building blocks of the Lox barcode cassette. Each block contains a code element and half a Lox site on each side.

Depending on the orientation of the Lox sites, there are four possible types of blocks (Fig. 2 D). Three colours have been used to code these: red, green and blue. By definition, the reverse complement of a block is of the same colour class. In contrast to blue blocks, red and green blocks have their Lox cores cleaved in a way such that their flanking Lox sites are unchanged after inversion, while the intervening sequence is reverse-complemented.

Blocks are similar to the concept of units in [27], introduced to derive expressions for the total number of sequences, stable or unstable, generated from a Lox cassette where all  $(n+1)$  sites can interact. Their analysis implies  $m > 82$  and a code diversity of order  $n$ . Quite distinctly, here we focus on enumerating size-stable sequences that arise in the regime  $5 \text{ bp} \leq m < 24 \text{ bp}$  with code diversities of order  $n^3$ .

Stable codes are necessarily made of blocks from the initial cassette, and as shown in Fig. 2 E, their composition in terms of block colors is prescribed. Letting  $n_r$ ,  $n_g$ , and  $n_b$  be the number of red, green, and blue blocks in the initial cassette with  $n$  elements, an upper bound on the number of possible barcodes of size  $k$  with  $k_r$  red,  $k_g$  green and  $k_b$  blue blocks is the number of possible outcomes when choosing  $k_r$ ,  $k_g$  and  $k_b$  from  $n_r$ ,  $n_g$  and  $n_b$  elements in arbitrary order:

$$k_r! \binom{n_r}{k_r} k_g! \binom{n_g}{k_g} k_b! \binom{n_b}{k_b} 2^{k_r+k_g},$$

where  $n_r + n_g + n_b = n$  and  $k_r + k_g + k_b = k$ . The additional factor  $2^{k_r+k_g}$  arises as there are two valid orientations of every code element of a red and green block after recombination. Conditioned on  $n_r$ ,  $n_g$ , and  $n_b$ , to derive an upper bound for a cassette's diversity, we add the numbers for the four possible stable barcode configurations of  $k_r$ ,  $k_g$ , and  $k_b$  (Fig. 2 E), taking into account that certain configurations appear more than once (e.g. the configurations with one red and two blue blocks appears three times). Using the expression above for each of the four configurations, for  $5 \text{ bp} \leq m < 24 \text{ bp}$ , and cassettes with inverted flanking sites pointing at each other (the opposite case is similar) this yields,

$$3 \left( 1! \binom{n_r}{1} 2! \binom{n_b}{2} \right) 2^1 + 1 \left( 2! \binom{n_r}{2} 1! \binom{n_g}{1} \right) 2^{2+1} + 2 \left( 1! \binom{n_r}{1} 1! \binom{n_b}{1} \right) 2^1 + 1 \left( 1! \binom{n_r}{1} \right) 2^1.$$

By construction,  $n_g = n_r - 1$ , and since  $n_b = n - 2n_r + 1$ , substituting the respective terms leads to an expression that is a function of  $n$  and  $n_r$  alone. For given  $n$  odd, this reduces the task of finding the optimal cassette design to an explicitly solvable one-dimensional optimization problem:

$$\arg \max_{n_r} \quad 32n_r^3 - 12(2n+3)n_r^2 + (6n^2 + 10n + 14)n_r \quad \text{for} \quad n_r \leq \frac{n+1}{2}.$$

For  $n \geq 5$ , the global maximum is achieved at the boundary  $n_r = (n+1)/2$ . This implies  $n_b = 0$ , and a global upper diversity bound of  $(n+1)(n-1)^2 + (n+1)$ , of order  $O(n^3)$ . It is easily verified that  $n_b = 0$  is only possible if the cassette design is alternating and  $n$  is odd, which implies the flanking sites are inverted.

## Alternating Lox cassette design achieves the upper diversity bound

For an alternating cassette design, achieving the code diversity upper bound requires complete freedom in code generation via recombination events. By construction, we show that this is the case if  $n \geq 7$ .

Consider an alternating cassette with five elements and  $m \geq 5 \text{ bp}$ , and recombination events that do not alter the size of the cassette (i.e., inversions). First note that red blocks in position three and five can move into the first position via a single recombination event. Furthermore, a red block in position one can be inverted by first moving to position three, then to five, and back again. A straight-forward recipe to create an arbitrary code made of a single red block is then to: i) move the block into the first position (if required); ii) change its orientation (if required); and finally iii) excise the remaining blocks.

Similarly, to generate an arbitrary code composed of a red and a green block from an alternating cassette with six elements, we can perform steps i) and ii). Then we apply the same procedure to the green blocks, leaving the first block untouched. This results in the first two blocks of the cassette being the desired code. To generate the size-stable code, elements that are not part of the code are excised.

Finally, for a cassette with seven elements, sequentially following the recipe given above, the first three blocks can be populated such that they match any possible code before excising the remaining blocks. This shows that any possible code of size one to three can be created via Lox recombination if the cassette is alternating,  $n \geq 7$ ,  $m \geq 5$  bp, and flanking sites are inverted.

Under constitutive Cre expression, barcodes with three elements can still undergo inversions via the flanking sites, which reduces their code diversity by a factor of two. The code diversity is therefore that given in Eq. (1).

## Design of code element sequences

That barcodes generated from a Lox cassette are pre-defined in terms of sequence and position in the genome represents an advantage over barcoding systems that rely on insertion site analysis for barcode readout [25, 40]. If codes-reading was error-free, choosing code elements of a particular color (Fig. 2 D) from a set of sequences that differ at least by one bp pair in both orientations would be sufficient. The maximum number of such elements is  $(4^m - 4^{m/2})/2$  and  $4^m/2$  for  $m$  even or odd, respectively, which is large even for small  $m$ .

In order to be perfectly robust to  $j$  read errors via nearest-neighbor match, all pairs of elements of a given color need to differ by a Hamming distance of at least  $2j + 1$  bp [41]. The size of the sets of elements that meet this condition quickly decreases with increasing  $j$  (see Fig. 3 A for numerical estimates). To ensure correction of two sequencing errors requires  $m \geq 5$  bp.

Assuming that sequencing errors arise independently and error rates are identical for all bases, the number of read sequencing errors in a code element of size  $m$  is Binomial with the error probability per bp [42]. Any element that has  $j$  or less errors will be classified correctly by nearest-neighbor matching. The probability of more than  $j$  errors gives an upper bound for the expected proportion of misclassified code elements. Fig. 3 B shows this for elements of size  $m = 7$  bp as a function of the minimal distance and the read error rates for next-generation sequencing platforms [43]. Different symbols indicate different sequence data. Even for low-fidelity platforms like Pacific Bioscience single molecule real time sequencing, a minimal distance of five bp results in less than ten misclassified elements per million.

## Probabilistic features of optimal Lox cassettes

In this section we explore stochastic features of the optimal design, specifically the probabilities to generate each of the final codes and the number of recombination events that are needed to create size-stable codes. For the analysis, we make two assumptions: first, all interactions with Lox sites that are at least 82 bp apart are equally likely; second, recombination events occur sequentially and independently.

### Barcode distribution is heterogeneous

Size-stable barcodes of a Lox cassette are randomly generated and not all codes are equally likely. Although an analytical expression for the probability mass function of final codes is not available, stochastic simulations enable us to study properties of practical importance such as the probability of generating a code more than once. Ensuring this probability is low is important in practice because progeny of two cells that independently generate the same code will be confounded as pertaining to the same clone.

Fig. 3 C shows the generation probability for each of the 1022 codes from a cassette with 13 elements. To produce this plot,  $10^8$  barcodes were Monte Carlo generated *in silico* via sequential recombination of the initial cassette. The number of times a specific code appeared was recorded, normalized and sorted. While some codes are relatively frequent, most are rare. In Fig. 3 D, the average number of recombination events (inversions: blue, excision: black) is plotted as a function of barcode probability. The number of inversions and barcode probability are negatively correlated, an indication that rare codes undergo, on average, more inversions. The number of excisions is close to two for all codes.

Ideally, each cell is tagged with a unique barcode. As with all existing barcoding techniques however, 100% unique barcodes cannot be guaranteed. What influences the expected number of unique barcodes is the code diversity  $D$ ,  $p_i$ , the probability of code  $i$ , where  $i \in \{1, 2, \dots, D\}$ , and  $j$ , the total number of codes that are generated. Using analysis of the generalized birthday party problem [44], the expected

proportion of unique codes is

$$\sum_{i=1}^D p_i (1 - p_i)^{j-1} \approx 1 - (j-1) \sum_{i=1}^D p_i^2, \quad (2)$$

where the numerically convenient approximation on the right hand side arises from a Taylor expansion around 0 and is appropriate if  $(j-1) \ll 1/(\max_i p_i)$ . Relatively large  $p_i$ 's negatively affect the expected proportion of unique codes. For heterogeneous barcode distributions, a natural strategy is to discard most frequent codes from the analysis. Barcodes that are included in the final analysis are called informative.

Using the approximation Eq. (2), in Fig. 3 E we computed the maximum number of cells that can initially be barcoded versus the number of cells that generate an informative code, for one to three sequential cassettes (indicated by the numbers 1, 2, 3), with the requirement that no more than 1% of informative codes are generated more than once. The color represents the percentage of discarded codes relative to the total code diversity. This parameter can be adjusted to meet the needs of a given experiment. E.g., for three concatenated cassettes with 13 elements each,  $10^5$  informative codes that are 99% unique can be generated by inducing barcodes in either  $10^6$  cells and including most codes or inducing barcodes in  $10^{12}$  cells and discarding most codes from the analysis. These results show that by discarding frequent codes from the read-out, large numbers of clones can be confidently tracked, indicating this *in situ* barcoding is suitable for high-throughput lineage tracing experiments.

### Number of recombination events to generate barcodes does not diverge with cassette size

If Cre is expressed for long enough, Lox cassettes will eventually become size-stable. The time this will take correlates with the number of recombination events that separate a stable barcode from its initial cassette. Below, we estimate this quantity using the theory of absorbing Markov chains.

In a cassette with  $n$  elements, there are  $n+1$  Lox sites. The number of Lox pairs that are flanking  $k$  elements is  $n+1-k$ . Lox pairs that have less than three elements between them do not interact as they are separated by less than the minimal distance. Pairs of Lox sites that have three or more elements between them are termed productive. For  $n \geq 3$  the number of productive pairs is  $\sum_{k=3}^n (n+1-k) = (n-1)(n-2)/2$ , and the number of productive pairs, where recombination leads to excision, i.e. where an even number of elements separates the two sites, is

$$\sum_{3 \leq k \leq n, k \text{ even}} (n+1-k) = \frac{(n-1)(n-3)}{4}$$

for  $n$  odd. The probability that a productive pair excises exactly  $k$  elements is given by the ratio of productive pairs that are separated by  $k$  elements to the total number of productive pairs, i.e.

$$P(\text{excision of } k \text{ elements}) = \frac{n+1-k}{\sum_{k=3}^n (n+1-k)} = \frac{2(n+1-k)}{(n-1)(n-2)}, \quad (3)$$

for  $k$  even,  $3 \leq k \leq n$ , otherwise it is zero. The number of productive pairs where recombination leads to inversion is (for  $n$  is odd)

$$\sum_{3 \leq k \leq n, k \text{ odd}} (n+1-k) = \frac{(n-1)^2}{4}, \quad (4)$$

and the probability that interaction of a productive pair leads to an inversion is

$$P(\text{inversion}) = \frac{2(n-1)^2}{4(n-1)(n-2)} = \frac{n-1}{2(n-2)}. \quad (5)$$

Equations (3) - (5) allow the formulation of size-stable barcodes as a discrete-time absorbing Markov chain. The number of elements in the cassette corresponds to its state, and Eq. (3) and Eq. (5) give the transition probabilities from  $n$  to  $n-k$ , and from  $n$  to  $n$  elements respectively. There are  $n-3$  transient and 4 absorbing states. Absorbing states are cassettes that have either three, two, one, or zero elements. Absorbing Markov models are well understood, and a wealth of theoretical predictions regarding their properties are available [45]. The fundamental matrix of this Markov Chain is

$$N = (I_{n-3} - Q)^{-1},$$

where  $I_{n-3}$  is an  $(n-3) \times (n-3)$  identity matrix, and  $Q$  is the transition matrix corresponding to the transient states. The expected number of recombination events, starting with a cassette of  $n$  elements, until reaching a final code is the  $n^{\text{th}}$  entry of the vector  $t = Nc$ , where  $c$  is a column vector all of whose entries are 1.

In Fig. 3 F, the average number of recombination events from the initial cassette to final code is shown as a function of the cassette length. Although code diversity grows as  $O(n^3)$ , the number of recombination events code generation increases linearly in  $n$ .

## Discussion

### Lox barcode cassettes with code elements of size four

When we identify the optimal Lox barcode cassette, we focus on code elements in the regime  $5 \text{ bp} \leq m < 24 \text{ bp}$ . These have maximal size-stable barcodes of three elements that are insensitive to over-estimation of the minimal Lox interaction distance. For  $m < 5 \text{ bp}$ , size-stable barcodes of four elements are possible and their maximal code diversity grows as  $O(n^4)$ . These are stable, however, only if the minimal interaction distance between two Lox sites is greater than 80 bp, a distance at which interactions have shown to still be possible *in vivo* in the similar Flp/FR system [31].

Most interesting is the case  $m = 4 \text{ bp}$ , which permits correction of one sequencing error with six code elements that are 3 bp apart in both orientations (see gray bars in Fig. 3 A). The upper diversity bound is derived along the same lines as for  $m \geq 5 \text{ bp}$  (see Fig. 4 E for possible stable codes), which gives

$$48 \binom{n_r}{1} \binom{n_b}{3} + 64 \binom{n_r}{2} \binom{n_g}{1} \binom{n_b}{1} + 12 \binom{n_r}{1} \binom{n_b}{2} + 16 \binom{n_r}{2} \binom{n_g}{1} + 4 \binom{n_r}{1} \binom{n_b}{1} + 2 \binom{n_r}{1}.$$

To maximize usage of the 6 code elements, we start with a cassette that has six red, five green and six blue blocks, i.e.  $\{n_r, n_g, n_b\} = \{6, 5, 6\}$ . This gives an upper diversity bound of 36996 barcodes. As confirmed by simulations, this upper bound is attained by a cassette with inverted flanking sites in which the first 11 Lox sites are alternating, and the remaining sites, except the last, are oriented in the same direction as the first Lox site (Fig. 4 F). Under constitutive Cre expression, barcodes with four elements can still undergo inversions, and the effective code diversity is 19,716.

Careful measurements will be needed to determine whether Lox sites at a distance of 80 bp still interact. If they don't, the cassette shown in Fig. 4 F with  $m = 4 \text{ bp}$  represents an interesting alternative to the design described in the main text, as with less elements it reaches higher code diversity, but at the cost of less robustness to sequencing error and hence barcode readout fidelity.

### Higher order Lox interactions

Single recombination events always involve exactly two Lox sites. However nothing except DNA flexibility prevents several pairs of Lox sites to interact simultaneously. The rate at which pairs of Lox sites bind depends on the number of Lox sites and the kinetic rates of Lox-Lox complexes. *In vitro*, the latter appear stable [32] and with the potentially large number of Lox sites in the barcode cassettes, make simultaneous interactions a plausible possibility.

Higher order Lox interactions lead to previously unreported, and in certain cases novel, recombination products (Fig. 4 C). For example, simultaneous interactions of two overlapping pairs of Lox sites oriented in the same direction do not result in excision, but in a reordering of the sequences between the sites. Similarly, if pairs are inverted, simultaneous recombinations do not invert but excise the sequence between the outermost sites.

For the alternating cassette and  $n \geq 7$ , multiple concurrent Lox interactions do not generate additional codes as the upper code diversity bound is already attained. Therefore our results on Lox barcode design and code elements remain unchanged in the presence of higher order Lox interactions. What changes is the distribution over barcodes, which flattens in the tail if more than one Lox pair recombines at a time (Fig. 4 D).

### Transient Cre expression.

Code diversity strongly depends on the number of elements in size-stable barcodes. If Cre is expressed constitutively, size-stable barcodes with code elements of size  $m \geq 5 \text{ bp}$  have a maximum of three elements. Another possibility is to create transient Cre activity rather than constitutive.



A well tested system that provides temporal control over Cre activity is tamoxifen inducible CreEr [34]. In the presence of tamoxifen, the fusion protein CreEr, which is normally located in the cytoplasm, is transported into the nucleus, where it can bind to Lox sites and induce recombination. Depending on the duration of Cre activation and its efficiency, stable sequences with more than three elements are likely to be generated from a Lox barcode cassette. Although most of these sequences are stable only in the absence of Cre, in this section we make no distinction between these and the size-stable barcodes defined earlier.

Fig. 4 A shows barcode probabilities after activation of CreEr in  $10^6$  cells with an optimal Lox cassette of size 13. The number of recombination events induced by transient CreEr activity is assumed Poisson with mean one. About  $10^4$  distinct barcodes are generated, and 30% of these appear only once. Although promising in terms of code diversity, it should be noted that potential drawbacks of this approach are the length of the barcodes (leading to more involved code sequencing), leakiness of CreEr into the nucleus in non-induced cells [46], and the relatively long half-life of tamoxifen [47].

Existing cellular barcoding approaches have already lead to significant biological discoveries and so new approaches that overcome their shortcomings are inherently desirable. Here we have established that using Cre Lox, it would be feasible to create an *in situ*, triggerable barcoding system with sufficient diversity to label a whole mouse, and propose this as a system for experimental implementation.

**Acknowledgments:** Ton Schumacher (Netherlands Cancer Institute) for informative discussions.

**Funding:** T.W., S.N and K.D. were supported by Human Frontier Science Program grant RGP0060/2012. K.D. was also supported by Science Foundation Ireland grant 12 IP 1263.

## References

- [1] Buchholz VR, Flossdorf M, Hensel I, Kretschmer L, Weissbrich B, Gräf P, et al. Disparate individual fates compose robust CD8+ T cell immunity. *Science*. 2013;340(6132):630–635.
- [2] Gerlach C, van Heijst JWJ, Swart E, Sie D, Armstrong N, Kerkhoven RM, et al. One naive T cell, multiple fates in CD8+ T cell differentiation. *J Exp Med*. 2010;207(6):1235–1246.
- [3] Verovskaya E, Broekhuis MJ, Zwart E, Ritsema M, van Os R, de Haan G, et al. Heterogeneity of young and aged murine hematopoietic stem cells revealed by quantitative clonal analysis using cellular barcoding. *Blood*. 2013;122(4):523–532.
- [4] Naik SH, Perié L, Swart E, Gerlach C, van Rooij N, de Boer RJ, et al. Diverse and heritable lineage imprinting of early haematopoietic progenitors. *Nature*. 2013;496(7444):229–232.
- [5] Perié L, Duffy KR, Kok L, de Boer RJ, Schmacher TN. The branching point in erythro-myeloid differentiation. *Cell*. 2015;163(7):1655–1662.
- [6] Johnson MB, Wang PP, Atabay KD, Murphy EA, Doan RN, Hecht JL, et al. Single-cell analysis reveals transcriptional heterogeneity of neural progenitors in human cortex. *Nat Neurosci*. 2015;18(5):637–646.
- [7] Yagi T. Genetic basis of neuronal individuality in the mammalian brain. *J Neurogenet*. 2013;27(3):97–105.
- [8] Zeisel A, Muñoz Manchado AB, Codeluppi S, Lönnerberg P, La Manno G, Juréus A, et al. Cell types in the mouse cortex and hippocampus revealed by single-cell RNA-seq. *Science*. 2015;347(6226):1138–1142.
- [9] Nolan-Stevaux O, Tedesco D, Ragan S, Makhanov M, Chenchik A, Ruefli-Brasse A, et al. Measurement of cancer cell growth heterogeneity through lentiviral barcoding identifies clonal dominance as a characteristic of in vivo Tumor engraftment. *PLoS One*. 2013;8(6):e67316+.
- [10] Bhang HeC, Ruddy DA, Krishnamurthy Radhakrishna V, Caushi JX, Zhao R, Hims MM, et al. Studying clonal dynamics in response to cancer therapy using high-complexity barcoding. *Nat Med*. 2015;21(5):440–448.
- [11] Klauke K, Broekhuis MJC, Weersing E, Dethmers-Ausema A, Ritsema M, González MV, et al. Tracing dynamics and clonal heterogeneity of Cbx7-induced leukemic stem cells by cellular barcoding. *Stem Cell Reports*. 2015;4(1):74–89.

- [12] Rohr JC, Gerlach C, Kok L, Schumacher TN. Single cell behavior in T cell differentiation. *Trends Immunol.* 2014;35(4):170–177.
- [13] Duffy KR, Hodgkin PD. Intracellular competition for fates in the immune system. *Trends Cell Biol.* 2012;22(9):457–464.
- [14] Hawkins ED, Markham JF, McGuinness LP, Hodgkin PD. A single-cell pedigree analysis of alternative stochastic lymphocyte fates. *Proc Natl Acad Sci USA.* 2009;106(32):13457–13462.
- [15] Rieger MA, Hoppe PS, Smejkal BM, Eitelhuber AC, Schroeder T. Hematopoietic cytokines can instruct lineage choice. *Science.* 2009;325(5937):217–218.
- [16] Gomes FL, Zhang G, Carbonell F, Correa JA, Harris WA, Simons BD, et al. Reconstruction of rat retinal progenitor cell lineages in vitro reveals a surprising degree of stochasticity in cell fate decisions. *Development.* 2011;138(2):227–235.
- [17] Giurumescu CA, Kang S, Planchon TA, Betzig E, Bloomekatz J, Yelon D, et al. Quantitative semi-automated analysis of morphogenesis with single-cell resolution in complex embryos. *Development.* 2012;139(22):4271–4279.
- [18] Duffy KR, Wellard CJ, Markham JF, Zhou JHS, Holmberg R, Hawkins ED, et al. Activation-induced B cell fates are selected by intracellular stochastic competition. *Science.* 2012;335(6066):338–341.
- [19] Richards JL, Zacharias AL, Walton T, Burdick JT, Murray JI. A quantitative model of normal *Caenorhabditis elegans* embryogenesis and its disruption after stress. *Dev Biol.* 2013;374(1):12–23.
- [20] Etzrodt M, Ende M, Schroeder T. Quantitative single-cell approaches to stem cell research. *Cell stem cell.* 2014;15(5):546–558.
- [21] Cohen AR. Extracting meaning from biological imaging data. *Mol Biol Cell.* 2014;25(22):3470–3473.
- [22] Lu R, Neff NF, Quake SR, Weissman IL. Tracking single hematopoietic stem cells in vivo using high-throughput sequencing in conjunction with viral genetic barcoding. *Nat Biotech.* 2011;29(10):928–933.
- [23] Naik SH, Schumacher TN, Perié L. Cellular barcoding: a technical appraisal. *Exp Hematol.* 2014;42(8):598–608.
- [24] Perié L, Hodgkin PD, Naik SH, Schumacher TN, de Boer RJ, Duffy KR. Determining lineage pathways from cellular barcoding experiments. *Cell Rep.* 2014;6(4):617–624.
- [25] Sun J, Ramos A, Chapman B, Johnnidis JB, Le L, Ho YJ, et al. Clonal dynamics of native haematopoiesis. *Nature.* 2014;514(7522):322–327.
- [26] Zador AM, Dubnau J, Oyibo HK, Zhan H, Cao G, Peikon ID. Sequencing the Connectome. *PLoS Biol.* 2012;10(10):e1001411+.
- [27] Wei Y, Koulakov AA. An exactly solvable model of random site-specific recombinations. *Bull Math Biol.* 2012;74(12):2897–2916.
- [28] Peikon ID, Gizatullina DI, Zador AM. In vivo generation of DNA sequence diversity for cellular barcoding. *Nucleic Acids Res.* 2014;42(16):e127.
- [29] Livet J, Weissman TA, Kang H, Draft RW, Lu J, Bennis RA, et al. Transgenic strategies for combinatorial expression of fluorescent proteins in the nervous system. *Nature.* 2007;450(7166):56–62.
- [30] Hoess R, Wierzbicki A, Abremski K. Formation of small circular DNA molecules via an in vitro site-specific recombination system. *Gene.* 1985;40(2-3):325–329.
- [31] Ringrose L, Chabanis S, Angrand PO, Woodroffe C, Stewart AF. Quantitative comparison of DNA looping in vitro and in vivo: chromatin increases effective DNA flexibility at short distances. *The EMBO Journal.* 1999;18(23):6630–6641.

- [32] Pinkney JN, Zawadzki P, Mazuryk J, Arciszewska LK, Sherratt DJ, Kapanidis AN. Capturing reaction paths and intermediates in Cre-loxP recombination using single-molecule fluorescence. *Proc Natl Acad Sci USA*. 2012;109(51):20871–20876.
- [33] Parrish M, Unruh J, Krumlauf R. BAC modification through serial or simultaneous use of CRE/Lox technology. *J Biomed Biotechnol*. 2011;2011:1–12.
- [34] Nagy A. Cre recombinase: the universal reagent for genome tailoring. *Genesis*. 2000;26:99–109.
- [35] Blattman JN, Antia R, Sourdive DJD, Wang X, Kaech SM, Murali-Krishna K, et al. Estimating the precursor frequency of naive antigen-specific CD8 T cells. *J Exp Med*. 2002;195(5):657–664.
- [36] Colvin GA, Lambert JF, Abedi M, Hsieh CC, Carlson JE, Stewart FM, et al. Murine marrow cellularity and the concept of stem cell competition: geographic and quantitative determinants in stem cell biology. *Leukemia*. 2004;18(3):575–583.
- [37] Sternberg N, Hamilton D, Hoess R. Bacteriophage P1 site-specific recombination. II. Recombination between loxP and the bacterial chromosome. *J Mol Biol*. 1981;150(4):487–507.
- [38] Hamilton DL, Abremski K. Site-specific recombination by the bacteriophage P1 lox-Cre system. Cre-mediated synapsis of two lox sites. *J Mol Biol*. 1984;178(2):481–486.
- [39] Guo F, Gopaul DN, Van Duyne GD. Structure of Cre recombinase complexed with DNA in a site-specific recombination synapse. *Nature*. 1997;389(6646):40–46.
- [40] Bystrykh LV, Verovskaya E, Zwart E, Broekhuis M, de Haan G. Counting stem cells: methodological constraints. *Nat Meth*. 2012;9(6):567–574.
- [41] Cover TM, Thomas JA. *Elements of Information Theory*. Wiley-Interscience; 1991.
- [42] Li H, Ruan J, Durbin R. Mapping short DNA sequencing reads and calling variants using mapping quality scores. *Genome Res*. 2008;18(11):1851–1858.
- [43] Ross MG, Russ C, Costello M, Hollinger A, Lennon NJ, Hegarty R, et al. Characterizing and measuring bias in sequence data. *Genome Biol*. 2013;14(5):R51+.
- [44] Koot MR, Mandjes M. The analysis of singletons in generalized birthday problems. *Probab Eng Inform Sc*. 2012;26(2):245–262.
- [45] Grinstead CM, Snell JL. *Introduction to Probability*. 2nd ed. American Mathematical Society; 1997.
- [46] Kretzschmar K, Watt FM. Lineage Tracing. *Cell*. 2012;148(1-2):33–45.
- [47] Reinert RB, Kantz J, Misfeldt AA, Poffenberger G, Gannon M, Brissova M, et al. Tamoxifen-induced Cre-loxP recombination is prolonged in pancreatic islets of adult mice. *PLoS One*. 2012;7(3):e33529+.

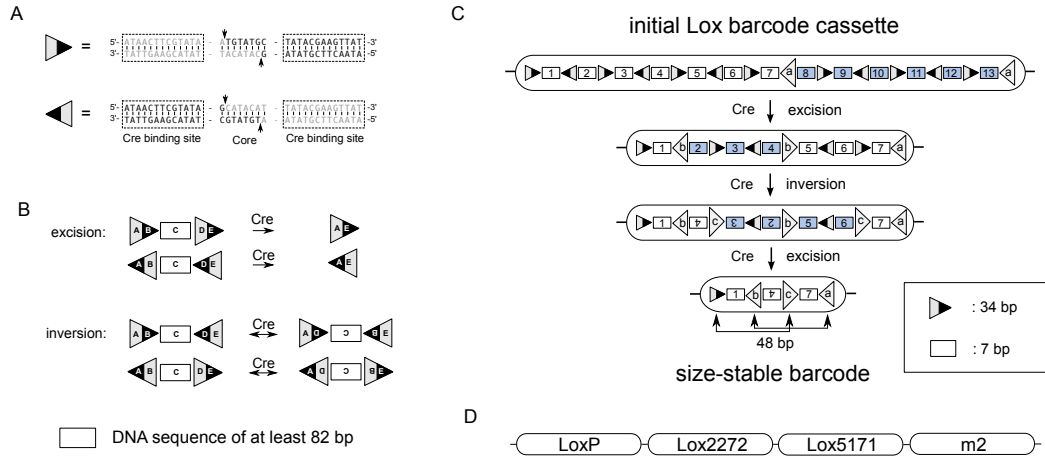


Figure 1: Lox biology and Lox barcode cassette. **A)** Lox DNA sequences. Lox sites are composed of two 13 bp palindromic Cre binding sites and an 8 bp core (original LoxP sequence shown). Cleavage sites in the core are indicated by arrows. **B)** Cre mediated site-specific excision and inversion of a sequence with a minimum of 82 bp between two Lox sites on the same chromosome [30]. If Lox sites are oriented in the same direction, recombination excises the sequence, while if they are oriented in opposite direction the sequence is inverted (i.e., the reverse complement). **C)** An alternating Lox cassette with 13 elements of size 7 bp. To illustrate how barcodes are generated, two excision and one inversion event are shown, creating a size-stable barcode with three random elements. Pairs of interacting Lox sites are indicated by a, b, and c. Elements affected by recombination have colored background. The barcode with three elements is size-stable as Lox sites oriented in the same direction (arrows) are closer than the minimal Lox interaction distance, precluding further excision. **D)** Four concatenated alternating Lox cassettes of 13 elements each with poorly-interacting Lox site variants [33] result in a code diversity greater than  $10^{12}$ .

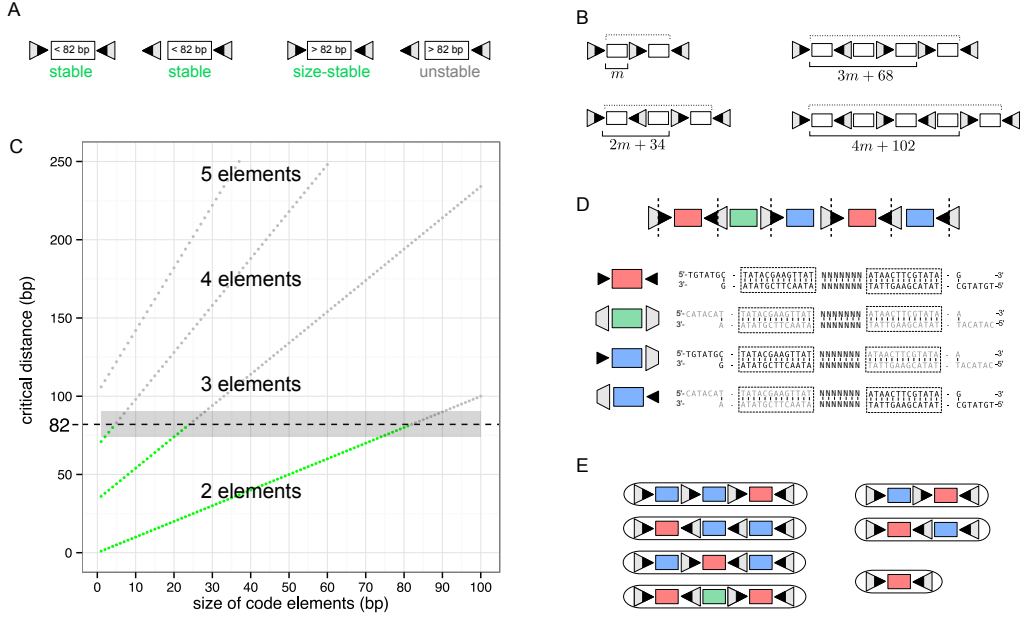


Figure 2: Barcode stability and code diversity. **A**) Size-stability of barcodes with a single element depends on the length of the sequence between the flanking sites and their relative orientation. **B**) Critical distances of barcodes of different sizes from a cassette with inverted flanking sites. Dotted lines show the critical distance if flanking sites are oriented in the same direction. **C**) Stability of barcodes from 2 to 5 elements for a Lox barcode cassette with inverted flanking sites. If the critical distance surpasses the minimal distance, stable codes (green) become unstable (gray). Barcodes of size three and four are unstable if  $m \geq 24$  and  $m \geq 5$  respectively, while codes of size five are always unstable. The gray interval illustrates potential uncertainty in the estimate of the minimal interaction distance. **D**) Sequences between Lox cleavage sites represent the fundamental building blocks of the barcode cassette. There are two with inverted Lox repeats (red, green) and two direct Lox repeats (blue) types of blocks. In the example, code elements are of size 7 bp and N denotes an arbitrary base. **E**) For a cassette with inverted flanking sites pointing at each other and  $5 \leq m < 24$ , four block compositions are possible ( $\{k_r, k_g, k_b\}$ ): two for barcodes of size three (three  $\{1, 0, 2\}$  and one  $\{2, 1, 0\}$ ), one for barcodes of size two (two  $\{1, 0, 1\}$ ) and one for barcodes of size one (one  $\{1, 0, 0\}$ ).

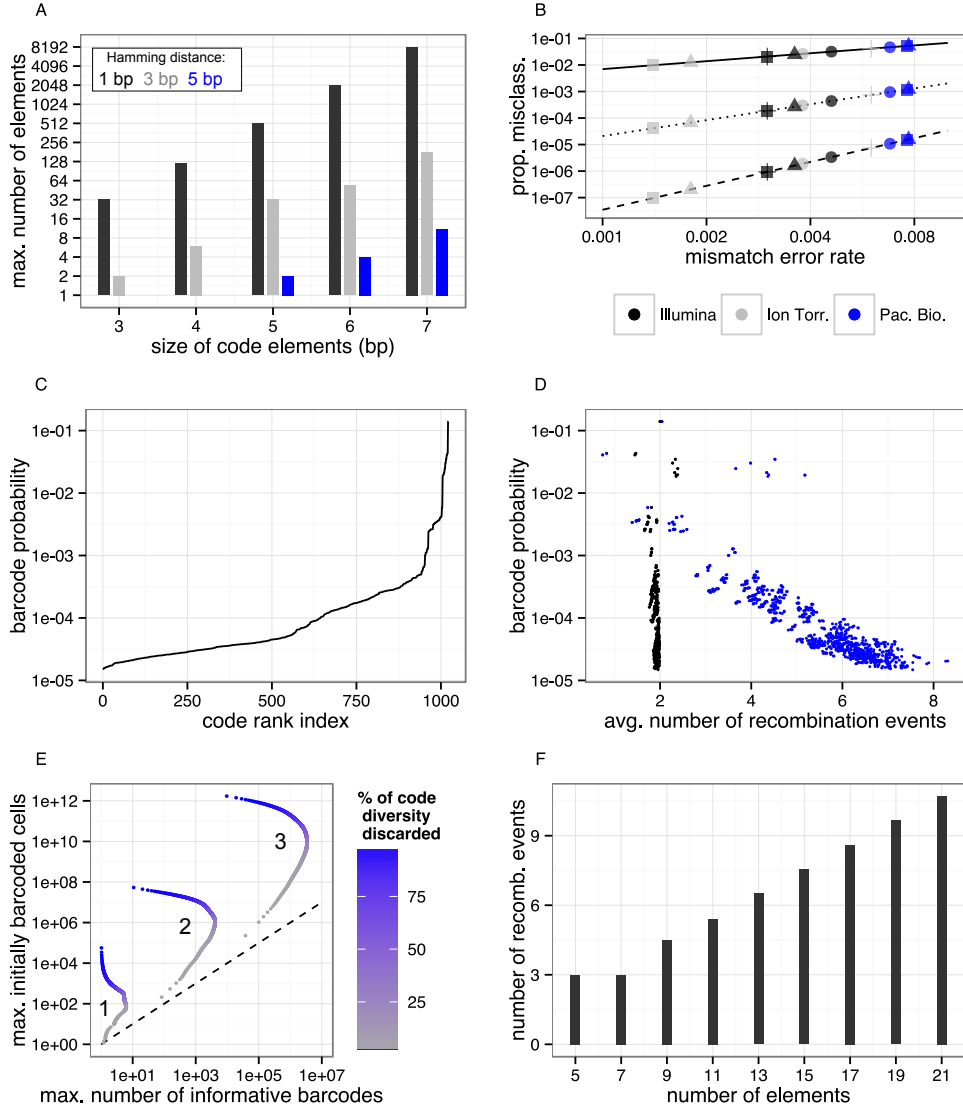


Figure 3: Design of code elements and probabilistic features of optimal cassettes. **A)** Computationally determined maximal size of sets of elements separated by a minimal Hamming distance of 1 bp (black), 3 bp (gray), and 5 bp (blue). To be robust to two sequencing errors, the minimal distance is 5 bp, which requires  $m > 4$  bp. **B)** Upper bound for the expected proportion of misclassified elements as a function of empirical DNA sequencing read error rates [43] for common sequencing platforms (Illumina, Ion Torrent, Pacific Biosciences) and different sequence data (*P. falciparum* (●), *E. coli* (▲), *R. spha.* (■), *H. sapiens* (+)). The minimal distance that separates the elements is 1 bp (solid), 3 bp (dotted), and 5 bp (dashed). **C)** Ranked probabilities of the 1022 size-stable barcodes from a cassette with 13 elements generated under constitutive Cre expression. A few codes are relatively frequent, but the majority are rare. **D)** Scatter-plot showing barcode probabilities against the average number of excisions (black) and the number of inversions (blue) that are generate size-stable barcodes from a 13 element optimal cassette. **E)** Number of cells in which a barcode can be induced versus the number of cells that produce informative codes, for one to three sequential cassettes, without exceeding 1% repeated occurrences in the informative codes. The color represents the percentage of discarded codes relative to the total code diversity, which can be adjusted to experimental conditions post acquisition. **F)** Although code diversity grows as  $O(n^3)$ , the expected number of recombination events that are needed to generate a size-stable code increases linearly in  $n$ .

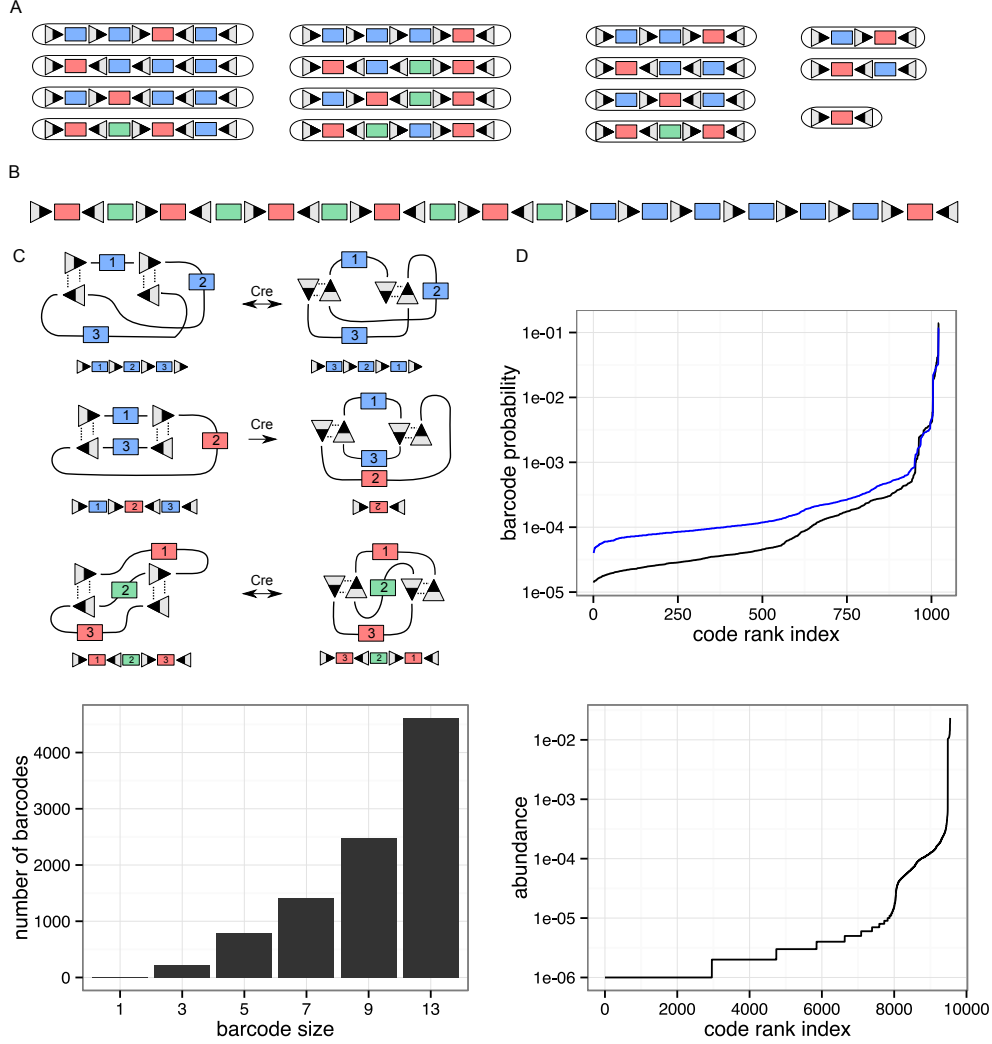


Figure 4: Short code elements, higher order Lox interactions, and transient Cre activation. **A)** Possible size-stable barcodes if  $m < 5$  bp. **B)** Cassette with 17 elements and  $m = 4$  bp that attains an effective code diversity of 19,716 barcodes if the minimal Lox interaction distance is greater than 80 bp. **C)** If two or more pairs of lox sites recombine simultaneously, unexpected recombination products can occur. **D)** Estimated barcode distribution if two Lox pairs can interact simultaneously (blue). The distribution becomes flatter at the lower end, implying that rare codes are more likely than if recombination events only occur sequentially (black). **E)** Mimicking a short Cre activation pulse in a population of a million cells carrying a 13 element Lox cassette, the number of recombination events is assumed Poisson distributed with mean 1. After the pulse many barcodes have not experienced any excisions. **F)** Code abundance after the pulse. Almost  $10^4$  distinct barcodes are generated, with 30% being generated once.