





# Boltzmann Sampling and Random Generation of Combinatorial Structures

Philippe Flajolet

Based on joint work with Philippe Duchon, Éric Fusy, Guy Louchard, Carine Pivoteau, Gilles Schaeffer

GASCOM'06, Dijon, September 12, 2006

C is a class of combinatorial structures.  $C_n$  = collection of objects of size n.

Draw uniformly at random from  $C_n$ ?:  $\mathbb{P}(\gamma) = \frac{1}{C_n}, \quad C_n := \|C_n\|.$ 

E.g.: trees, permutations, words, graphs, mappings, maps, etc.



Classification theory [Van Cutsem]; image synthesis [Viennot]; random testing in software eng. [J. Fayolle], combinatorics; simulation & statistical analysis of models in genetics [Denise], ecology [de Reffie], ...

Bijective method Surjective method Rejection method Markov method Recursive method

### Random Generation and Combinatorics

- *Bijective method*: find bijection with simpler (product) set.
- Surjective method: find a "multiple" set that is simpler
- Rejection method: find a larger set and filter.
- Markov method: superimpose Markov chain structure & travel!
- Recursive method: decompose according to counting probabilities
- Boltzmann: This talk!

Bijective method Surjective method Rejection method Markov method Recursive method

### **Bijective** method

Find bijection with simpler set Class C is such that  $C_n = \|C_n\|$  is a product.

**Words:**  $\mathcal{W}_n \cong \{a, b\}^n \Longrightarrow n$  random flips. **Permutations:**  $\mathcal{P}_n \cong [0] \times [0 \dots 1] \times \dots \times [0 \dots n-1] \Longrightarrow n$  RVs





Bijective method Surjective method Rejection method Markov method Recursive method

## Surjective method

Find many-to-one uniform correspondence between  $C_n$  and simpler set  $A_n$ .

divisibility:  $C_n \mid A_n$ .

Dyck excursions: by *conjugacy* with bridges  $\rightsquigarrow$  Catalan trees.

$$C_n=\frac{1}{2n+1}\binom{2n+1}{n}.$$

Jean-Luc Rémy's algorithm for binary trees.

Planar maps: cf Schaeffer et al.: by tree conjugation.

Usually requires pure product form!

Bijective method Surjective method Rejection method Markov method Recursive method

# Rejection method

Find larger set such that  $\mathcal{C}_n \subset \mathcal{D}_n$ , with simpler  $\mathcal{D}$ 

- $\implies$  Draw  $\delta \in \mathcal{D}$ . Test whether  $\delta \in \mathcal{C}$ ; repeat if needed Problem: Probability of success is  $\frac{C_n}{D_n}$ .
- E.g. Prime numbers; irreducible polynomials. Cf Ruskey.
- E.g. Florentine algorithm for Dyck/Motzkin meanders.



Bijective method Surjective method Rejection method Markov method Recursive method

## Markov method

- <u>View elements of a class  $S_n$  as states of a Markov chain</u>

— Set up transitions (e.g, via transformations)

If the graph is regular, then the stationary distribution is uniform.





Reversible Markov chains, Coupling [Propp-Wilson, Jerrum,...]. ~ Self-avoiding walks, dimer coverings, "hard" combinatorial objects.

May need information on mixing speed  $\lambda_2$ .

**Bijective method** Surjective method Rejection method Markov method Recursive method

## Recursive method

- Use counting sequences to decide splitting probabilities.
- E.g.: **Binary trees** with *n* external nodes, class  $\mathcal{B}_n$ .
- A. Set up recurrence  $B_n = \sum_{k=1}^{n-1} B_k B_{n-k}$ . B. Split  $n \mapsto \langle k, n-1-k \rangle$  with probability  $\frac{B_k B_{n-k}}{B_n}$ .

### Theorem (Recursive method)

Complexity of preprocessing is  $O(n^2)$  large integer operations. Complexity of boustrophedonic random generation is  $O(n \log n)$ arithmetic operations.

- ECO systems. Wilf's path approach.
- J. van der Hoeven: Preprocessing in time  $O(n^{1+\varepsilon})$ . A. Denise &
- P. Zimmermann: Floating point implementations. Also: MAPLE Combstruct.

### Boltzmann framework

### **Principle:**

- Generate according to a **distribution spread over all** C, depending on **control parameter** x.
- Size becomes a random variable (RV).
- Target choice of x to get objects of size near n with fair probability.

Cf Statistical Physics:  $\mathbb{P}(\gamma) = \frac{1}{Z} \exp\left(-\frac{\beta}{T} E[\gamma]\right).$ 

Ordinary (unlabelled) Boltzmann models

Assign to  $\gamma \in C$  probability proportional to exponential of its size:

$$\mathbb{P}(\gamma) \propto x^{|\gamma|} \implies \mathbb{P}(\gamma) = rac{x^{|\gamma|}}{C(x)},$$

 $C(x) = \sum_{n} C_{n} x^{n}$  is ordinary generating function (OGF). Requires  $x \le \rho_{C}$ , where  $\rho_{C}$  is the radius of convergence of C(x).

→ Size becomes a random variable:

$$\mathbb{P}(\mathsf{Size}=n)=\frac{C_n x^n}{C(x)}.$$

Unions, products, and sequences Labelled models, sets and cycles Unlabelled sets and cycles

## Boltzmann Samplers: the Plan!

Develop design rules given combinatorial specifications.

- Basic constructions:  $\cup, \times, SEQ$
- Labelled models: add SET, CYC
- Return to unlabelled models: add MSET, PSET, CYC

Do optimization w.r.t. size at the end: complexity issues.

Based on [DuFILoSc04] in CPC for labelled; [FIFuPi06] for unlabelled. Cf. F.+Sedgewick, *Analytic Combinatorics*.

Unions, products, and sequences Labelled models, sets and cycles Unlabelled sets and cycles

## Unions, products

### Lemma (Disjoint unions)

Boltzmann sampler  $\Gamma C$  for  $C = A \cup B$ : With probability  $\frac{A(x)}{C(x)}$  do  $\Gamma A(x)$  else do  $\Gamma B(x)$ 

### Lemma (Products)

Boltzmann sampler  $\Gamma C$  for  $C = \mathcal{A} \times \mathcal{B}$ : Generate independent pair  $\langle \Gamma \mathcal{A}(x), \Gamma \mathcal{B}(x) \rangle$ .

*Proofs* = *One-liners!* Using basic definitions of probability.

— Disjoint union:  $|\gamma| = n \Longrightarrow$  if  $\gamma \in \mathcal{A}$  then  $\mathbb{P}_{\mathcal{C}}(\gamma) = \frac{x^n}{A(x)} \cdot \frac{A(x)}{C(x)} \dots$ 

- Product: 
$$\mathbb{P}_{\mathcal{C}}(\gamma) = \frac{x^n}{A(x)} \cdot \frac{x^{n-n}}{B(x)} = \frac{x^n}{C(x)}$$

Unions, products, and sequences Labelled models, sets and cycles Unlabelled sets and cycles

## Sequences

### Lemma (Sequences)

Boltzmann sampler  $\Gamma C$  for C = Seq(A):

- Generate K which is geometric with parameter A(x)
- Generate independent K-tuple  $\langle \Gamma A(x), \ldots, \Gamma A(x) \rangle$ .

**Proof.** Recursive equation:  $C = \mathbf{1} + AC$  with  $+, \times$  constructions. With probability  $\frac{1}{A(x)}$  STOP; else  $\Gamma A(x)$  and continue rec. with  $\Gamma C(x)$ . Number of trials of Bernoulli RV till success is Geometric.

Unions, products, and sequences Labelled models, sets and cycles Unlabelled sets and cycles

# Specifications with $\{\cup, \times, SEQ\}$

Specs	GF	Sampler
$1$ or $\mathcal{Z}$ (atom)	1 or <i>x</i>	$\Gamma \mathcal{C}:=$ output $f 1$ or $ullet$
$\mathcal{C}=\mathcal{A}\cup\mathcal{B}$	C(x) = A(x) + B(x)	$\Gamma C(x) := \frac{A(x)}{C(x)} \longrightarrow \Gamma B(x) \mid \Gamma C(x)$
$\mathcal{C}=\mathcal{A} imes\mathcal{B}$	$C(x) = A(x) \times B(x)$	$\Gamma C(x) := \langle \Gamma B(x), \Gamma C(x) \rangle$
$\mathcal{C}=\operatorname{Seq}(\mathcal{A})$	$C(x) = \frac{1}{1 - A(x)}$	$\Gamma C(x) := \operatorname{Geom}[A(x)] \Longrightarrow \Gamma A(x)$

Compile sampler from specification automatically.

Unions, products, and sequences Labelled models, sets and cycles Unlabelled sets and cycles

# Specifications with $\{\cup, \times, SEQ\}$ — continued

### Theorem (Complexity Minitheorem)

*Given* **oracle** *that provide the finitely many values of GFs, complexity is* **linear** *in size of object produced.* 

**Proof**  $\{\cup, \times, SEQ\}$ : overhead O(1) per node of derivation tree. Complexity model: exact computations over  $\mathbb{R}$ ; in practice, "floats" (more later).

#### Definition

 $\begin{array}{l} \mbox{Regular specification} = \mbox{iterative (nonrecursive) with } \{\cup,\times,{\rm SEQ}\}.\\ \mbox{Contex-free specification} = \mbox{recursive with } \{\cup,\times,{\rm SEQ}\}. \end{array}$ 

### Proposition

**Regular structures** and **context-free structures** have Boltzmann samplers of linear-time complexity.

Unions, products, and sequences Labelled models, sets and cycles Unlabelled sets and cycles

Specifications with  $\{\cup, \times, SEQ\}$  — continued (2)

### **Regular specifications**

• **Binary words** with **longest run** of a's of length < 17.

 $\operatorname{SeQ}_{<17}(\{a\}) \cdot \operatorname{SeQ}(b\operatorname{SeQ}_{<17}(\{a\})).$ 

- Codes, e.g., {*aba*, *abaaa*, *abba*}.
- Polyominos that have rational GF, e.g., Vertically convex.
- Languages recognized by deterministic finite automata E.g., Strings containing three times the pattern "abracadabra".
- Paths in digraphs even in the presence of sinks.

Unions, products, and sequences Labelled models, sets and cycles Unlabelled sets and cycles

Specifications with  $\{\cup, \times, SEQ\}$  — continued (3)

### **Contex-free specifications.**

- Binary trees:  $\mathcal{B} = \mathcal{Z} + \mathcal{B} \times \mathcal{B}$ .
- Solve quadratic equation  $B = x + B^2$  numerically, given x;
- Out put single node with probability  $\frac{x}{B}$ ;
- Else: Do two independent recursive calls to  $\Gamma B(x)$ .

For **rooted unlabelled trees**, Boltzammn model reduces to branching process.

Generate Motzkin trees [≠Alonso-Schoot], (unbalanced)

2-3-trees; random walks with finite step sets (dice), etc.

**Noncrossing graphs:** 

Unions, products, and sequences Labelled models, sets and cycles Unlabelled sets and cycles

Exponential (labelled) Boltzmann models

• For labelled classes, model is called exponential or labelled Boltzmann model

$$\mathbb{P}(\gamma) \propto \frac{x^{|\gamma|}}{|\gamma|!} \implies \mathbb{P}(\gamma) = \frac{1}{C(x)} \frac{x^{|\gamma|}}{|\gamma|!}$$
$$C(x) := \sum_{n} C_{n} \frac{x^{n}}{n!} \text{ is exponential GF (EGF).}$$

- Replace Cartesian product by labelled product (distribute labels).
- Unions, products, sequences: work like before, but with EGFs.
- Sets and cycles = to do!

Unions, products, and sequences Labelled models, sets and cycles Unlabelled sets and cycles

### Labelled sets and cycles

Poisson law: 
$$\mathbb{P}(X = k) = e^{-\lambda} \frac{\lambda^k}{k}!$$
.  
Logarithmic law:  $\mathbb{P}(X = k) = \frac{1}{l} \frac{\lambda^k}{k}$ ,  $L := 1/\log(1 - \lambda)^{-1}$ .

#### Lemma

Labelled sets and labelled cycles are obtained by a Poisson and Logarithmic generator resp.

 $\begin{array}{rcl} \mathcal{C} = \operatorname{Set}(\mathcal{A}) & : & \operatorname{Pois}(\mathcal{A}(x)) \Longrightarrow \mathsf{\Gamma}\mathcal{A}(x) \\ \mathcal{C} = \operatorname{Cyc}(\mathcal{A}) & : & \operatorname{Loga}(\mathcal{A}(x)) \Longrightarrow \mathsf{\Gamma}\mathcal{A}(x) \end{array}$ 

Cf: C = SEQ(A) :  $Geom(A(x)) \Longrightarrow \Gamma A(x)$ .

Applies to any specifiable class of combinatorial objects

- For each x, need finite # of computable real constants.
- Linear-time random generation.
- Size is not controlled (yet)
- Example: Cayley trees =  $\mathcal{T} = \mathcal{Z} \star \text{Set}(\mathcal{T})$ .
- Solve  $T(x) = xe^{T(x)}$  numerically.
- Generate root  $(\mathcal{Z})$ ;
- Choose random root degree as  $\Delta := \text{Pois}(T(x))$ ;
- Call  $\Delta$  independent copies of  $\Gamma(x)$ ;
- Hope for the best regarding size ( $\rightsquigarrow$  later)

Unions, products, and sequences Labelled models, sets and cycles Unlabelled sets and cycles

Examples:

Set partitions.  $S = SET(SET_{\geq 1}(Z))$ . # components is  $Pois(e^x - 1)$ ; each comp. is  $Pois(x) \mid \geq 1 = Vershik$ .

Ordered set partitions. Geometric triggers Poisson.

Assemblies of filaments. Poisson triggers geometric.

Unions, products, and sequences Labelled models, sets and cycles Unlabelled sets and cycles

### Unlabelled sets and cycles



[Pólya] Carbon has valency 4; hydrogen has valency 1. How to generate a **random alcohol**?.

= Nonplane unlabelled tree with node degrees  $\in \{0,3\}$ . Need to take care of symmetries to generate object only once!

Unions, products, and sequences Labelled models, sets and cycles Unlabelled sets and cycles

• The multiset construction C = MSET(A): form all finite multisets,

 $\mathcal{C} \cong \prod_{\alpha \in \mathcal{A}} \operatorname{Seq}(\{\alpha\}).$ 

(i) Gedanken Alg. Scan A & generate α with multiplicity Geom(x<sup>|α|</sup>).
(ii) Observe GF equation: C(x) = exp(A(x)) · exp(<sup>1</sup>/<sub>2</sub>A(x<sup>2</sup>)) · · · .
(iii) Do Poisson-controled generator for A with parameter A(x); repeat with <sup>1</sup>/<sub>2</sub>A(x<sup>2</sup>); etc.
(iv) Compute when to stop. Collect multset.

Proof involves  $\text{Geom}(\lambda) \equiv \text{Pois}(\lambda) + \text{Pois}(\frac{1}{2}\lambda^2) + \cdots$ .

Unions, products, and sequences Labelled models, sets and cycles Unlabelled sets and cycles

### Powersets and cycles

• The cycle construction: proceed from GFs. For C = CYC(A),

$$C(z) = \log \frac{1}{1 - A(z)} + \frac{1}{2} \log \frac{1}{1 - A(z^2)} + \cdots$$

Treat as infinite union, cf multisets. E.g., Necklaces.

- The **powerset** construction C = PSET(A): form all finite sets (no repetition!). Use identity  $1 + z = \frac{(1-z^2)}{(1-z)}$ . Generate Boltzmann multiset and throw away all elements of even multiplicity.
- Relativized constructions like  $C = MSET_3(A)$ : do  $\Gamma A(x^3)$ , etc.

Unions, products, and sequences Labelled models, sets and cycles Unlabelled sets and cycles

## Unlabelled constructions

### Theorem (Main Complexity Theorem)

For a class C specified (poss. recursively) from finite sets using

 $+, \times, \text{SEQ}, \text{MSET}, \text{MSET}_k, \text{CYC}, \text{CYC}_k,$ 

The Boltzman sampler  $\Gamma C(x)$  operates in linear time in the size of the object produced.

Also allow for powersets as soon as  $\rho < 1$ . Examples. Integer partitions, nonplane unlabelled trees, alcohols, mapping patterns [functional graphs], series-parallel circuits, etc

Unions, products, and sequences Labelled models, sets and cycles Unlabelled sets and cycles



Unions, products, and sequences Labelled models, sets and cycles Unlabelled sets and cycles



Unions, products, and sequences Labelled models, sets and cycles Unlabelled sets and cycles



Size control Discrete samplers

## Complexity

• Size control

$$\mathsf{PGF}(\mathsf{Size}) = \frac{C(ux)}{C(x)} \implies \mathbb{E}_x(\mathsf{Size}) = \frac{xC'(x)}{C(x)}.$$

Usually requires  $x \rightarrow \rho_C$  to get large structures.

Size control Discrete samplers



Free Boltzmann samplers: produce objects with randomly varying sizes!
 E.g., VC-polyominos: 37, 158, 389, 91, 21, 110, ...

Size control Discrete samplers

# Size control (1)

- Free Boltzmann samplers: produce objects with randomly varying sizes! E.g., VC-polyominos: 37, 158, 389, 91, 21, 110, ...
- **Tuned Boltzmann samplers**: choose *x* so that expected size = *n*.

Size control Discrete samplers

# Size control (1)

- Free Boltzmann samplers: produce objects with randomly varying sizes! E.g., VC-polyominos: 37, 158, 389, 91, 21, 110, ...
- **Tuned Boltzmann samplers**: choose *x* so that expected size = *n*.
- Analysis of size distribution of free sampler determines complexity.

Size control Discrete samplers

# Size control (2)

### "Frequent" profiles: [cf Analytic Combinatorics]



Depends on *singularity type* of generating function.

Size control Discrete samplers

### Theorem (Complexity I)

"Bumpy type" is granted for Hayman-admissible models. Approximate-size complexity = O(n). Exact size =  $o(n^2)$ .

Applies to GFs that are of type  $Exp \circ Fast-growth$ .

#### Theorem (Complexity II)

"Flat type" is granted for algebraic-logarithmic sing. + infinite Approximate-size complexity = O(n). Exact-size =  $o(n^2)$ .

### Theorem (Complexity III)

For "critical sequences":

Exact-size complexity = O(n).

Renewal type of algorithm at critical  $\rho$ .

Size control Discrete samplers

### Size control (3): Pointing

**Pointing:** If  $\mathcal{A}$  is a class, then  $\mathcal{C} = \mathcal{A}^{\bullet}$  is the set of objects with one atom pointed, and

$$C_n = nA_n,$$
  $C(z) = z \frac{d}{dz}A(z).$ 

Uniformity at given size is preserved (only size profile is altered). Transforms peaked (inefficient) distributions to flat (efficient). E.g., **binary trees** B:

$$\mathcal{B} = \mathcal{Z} + \mathcal{B} imes \mathcal{B} \implies \begin{cases} \mathcal{B} = \mathcal{Z} + \mathcal{B} imes \mathcal{B} \\ \mathcal{B}^{ullet} = \mathcal{Z} + \mathcal{B}^{ullet} imes \mathcal{B} + \mathcal{B} imes \mathcal{B}^{ullet}. \end{cases}$$

All simple families of trees: it works!

Size control Discrete samplers

### **Discrete samplers**

• Real arithmetics versus bit [boolean] complexity?

— Do bit-level generators for Bernoulli, Geometric, Poisson, Logarithmic.

$$\frac{1}{\pi} = \langle 0.010100010111110011000 \rangle_2.$$

Bernoulli: return bit at position  $\text{Geom}(\frac{1}{2})$ ; Geometric: iterate till 1. Cf. Knuth-Yao (1976); Von Neumann. Soria-Pelletier et al.

- Integrated samplers for set partitions, etc? Expect low bit-complexity!
- In practice do 40D evaluations of constants and be happy!

Size control Discrete samplers

## Conclusions



A plane parttion of size 15,000 [Carine Pivoteau]

Size control Discrete samplers

## Some literature (all on the web!)

[DuFILoSc04] "Boltzmann Samplers for the Random Generation of Combinatorial Structures", by Philippe Duchon, Philippe Flajolet, Guy Louchard, Gilles Schaeffer. In *Combinatorics, Probability, and Computing,* Special issue on Analysis of Algorithms, 2004, Vol. 13, No 4–5, pp. 577-625. [FIFuPi06] "Boltzmann Sampling of Unlabelled Structures", by Philippe Flajolet, Éric Fusy, and Carine Pivoteau. 14 pages. Submitted to *ANALCO'07*. [BaNi06] "Accessible and deterministic automata: enumeration and Boltzmann samplers". F. Bassino C. Nicaud. In *Fourth Colloquium on Mathematics and Computer Science*.

[Fusy05] "Quadratic exact-size and linear approximate-size random sampling of planar graphs", by Éric Fusy. In *2005 International Conference on Analysis of Algorithms. DMTCS* Conference Volume AD (2005), pp. 125-138.

[BoFuPi06] "Random sampling of plane partitions". By Olivier Bodini, Éric Fusy, and Carine Pivoteau. Gln *ASCOM-2006*.