## STOCHASTIC CHEMICAL KINETICS

**Dan Gillespie**
GillespieDT@mailaps.org

*Current Support*:  **Caltech** (NIGMS)
                  **Caltech** (NIH)
                  **University of California at Santa Barbara** (NIH)

*Past Support*:  **Caltech** (DARPA/AFOSR, Beckman/BNCM))
                  **University of California at Santa Barbara** (DOE)
                  **Molecular Sciences Institute** (Sandia/DOE)
                  **Office of Naval Research**

*Main Collaborators*:
                  **Linda Petzold** (UCSB)
                  **John Doyle** (Caltech)
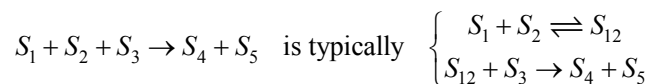
---

**A Chemically Reacting System** consists of …

- Molecules of $N$ **chemical species** $S_1, \ldots, S_N$.

  - Inside a volume $\Omega$, at some temperature $T$.

- $M$ "elemental" **reaction channels** $R_1, \ldots, R_M$.

  - We assume $R_j$ to be a ***single instantaneous physical event*** that changes the population of at least one species.

  - In practice, "elemental" means that $R_j$ must be one of two types:

          *Unimolecular*:    $S_i \to$ something else,
                    or
        *Bimolecular*:  $S_i + S_{i'} \to$ something else.

  - All other types of reaction (trimolecular, reversible, etc.) are made up of a series of two or more elemental reactions. E.g.:

$$S_1 + S_2 + S_3 \to S_4 + S_5 \quad \text{is typically} \quad \begin{cases} S_1 + S_2 \rightleftharpoons S_{12} \\ S_{12} + S_3 \to S_4 + S_5 \end{cases}$$

*Question*:  **How does a spatially homogeneous (or well-stirred) chemically reacting system evolve in time?**

*The Trad Answer*:
> **According to the *reaction rate equation* (RRE).**

- A set of coupled, first-order ODEs.

- Derived using ad hoc, phenomenological reasoning.
  - Is *more* than the "mass action equations" of thermodynamics, which apply only to systems in complete equilibrium.

- Implies the system evolves *continuously* and *deterministically*.
  - Yet molecules come in integer numbers and react stochastically.

- Is empirically accurate for large (test tube size) systems

- But is sometimes not adequate for very small (cell-size) systems.

* * *

**Doing it "right":  Molecular Dynamics**

- The most exact way of describing the system's evolution.
- The "motion picture" approach:  Tracks the position and velocity of every molecule in the system.
- Simulates *every* collision, *non-reactive* as well as *reactive*.
- Shows changes in species populations and their spatial distributions.
- ***But*** . . . it's *unfeasibly slow* for nearly all realistic systems.

**A great simplification occurs *if* successive *reactive* collisions tend to be separated in time by *very many non-reactive* collisions.**

- The overall effect of the non-reactive collisions is to *randomize*
  - the *velocities* of the molecules (Maxwell-Boltzmann distribution).
  - the *positions* of the molecules (spatially uniform or **well-stirred**),
- Then, instead of having to describe the system's state as the *position, velocity and species of each molecule*, we need only give

$$\mathbf{X}(t) \triangleq \left( X_1(t), \ldots, X_N(t) \right),$$

$$X_i(t) \triangleq \text{ the } \textit{number} \text{ of } S_i \text{ molecules at time } t.$$

---

But this ***well-stirred simplification***, which . . .

- *ignores* the non-reactive collisions,
- *drastically truncates* the definition of the system's state,

. . . comes at a price:

**$\mathbf{X}(t)$ *must now be viewed as a stochastic process*.**

➢ But in fact, ***the system was never deterministic to begin with.***
   Even if molecules moved according to classical mechanics . . .
   - Unimolecular reactions always involve randomness (QM).
   - Bimolecular reactions usually do too.
   - A system of many colliding molecules is so *sensitive to initial conditions* that, for all practical purposes, it evolves "randomly".
   - The system is not isolated. It's in a *heat bath*, which keeps it "at temperature $T$ " – via essentially random interactions.

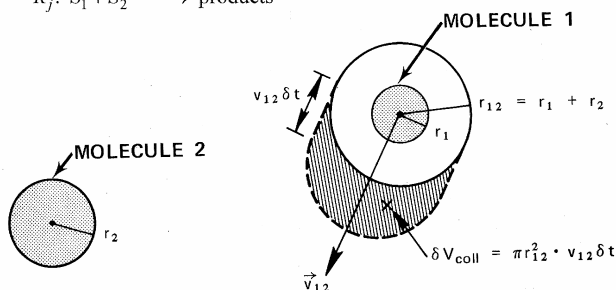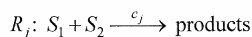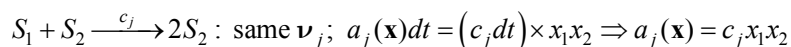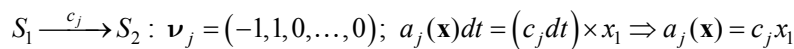**For well-stirred systems**, each $R_j$ is completely characterized by …

- a ***propensity function*** $a_j(\mathbf{x})$ : Given the system is in state $\mathbf{x}$,

  $a_j(\mathbf{x})\,dt \triangleq$ *probability* that one $R_j$ event will occur in the next $dt$.

  - The *existence* and *form* of $a_j(\mathbf{x})$ follow from *molecular physics*.

- a ***state change vector*** $\boldsymbol{\nu}_j \equiv \left(\nu_{1j},\ldots,\nu_{Nj}\right)$:

$$\nu_{ij} \triangleq \text{the } \textit{change} \text{ in } X_i \text{ caused by one } R_j \text{ event.}$$

  - $R_j$ induces $\mathbf{x} \to \mathbf{x} + \boldsymbol{\nu}_j$. $\left\{\nu_{ij}\right\} \equiv$ the "stoichiometric matrix."

***Examples*:**

$$S_1 \xrightarrow{c_j} S_2 : \; \boldsymbol{\nu}_j = \left(-1,1,0,\ldots,0\right); \; a_j(\mathbf{x})dt = \left(c_j dt\right) \times x_1 \Rightarrow a_j(\mathbf{x}) = c_j x_1$$

$$S_1 + S_2 \xrightarrow{c_j} 2S_2 : \; \text{same } \boldsymbol{\nu}_j; \; a_j(\mathbf{x})dt = \left(c_j dt\right) \times x_1 x_2 \Rightarrow a_j(\mathbf{x}) = c_j x_1 x_2$$

---

$R_j: \; S_1 + S_2 \xrightarrow{c_j} \text{products}$



$$\text{Prob}\{\text{collision in } dt\} = \frac{(\pi r_{12}^2)(\overline{v}_{12}dt)}{\Omega}. \quad \text{Prob}\{R_j \,|\, \text{collision}\} \equiv p_j.$$

$$a_j(\mathbf{x})\,dt \equiv \underbrace{\left(\frac{(\pi r_{12}^2)(\overline{v}_{12}dt)}{\Omega}\right) \times p_j}_{\substack{\text{Prob that a typical } S_1\text{-}S_2 \\ \text{pair reacts in next } dt}} \times x_1 x_2 = \underbrace{\left(\frac{\pi r_{12}^2 \overline{v}_{12} p_j}{\Omega}\right)}_{c_j} x_1 x_2 \, dt = \underbrace{c_j x_1 x_2 \, dt}_{a_j(\mathbf{x})}$$

$$R_j \text{ iff } \textit{"collisional K.E."} > E_j \; \Rightarrow \; p_j = \exp\left(-\frac{E_j}{k_B T}\right) \text{ … Arrhenius!}$$

***Diffusional motion (well-stirred)*:** $\; \pi r_{12}^2 \overline{v}_{12}$ is replaced by $4\pi r_{12}(D_1 + D_2)$

**Two exact, rigorously derivable consequences . . .**

> **1.** The *chemical master equation* (CME):

$$\frac{\partial P(\mathbf{x},t\,|\,\mathbf{x}_0,t_0)}{\partial t} = \sum_{j=1}^{M}\Big[a_j(\mathbf{x}-\boldsymbol{\nu}_j)P(\mathbf{x}-\boldsymbol{\nu}_j,t\,|\,\mathbf{x}_0,t_0) - a_j(\mathbf{x})P(\mathbf{x},t\,|\,\mathbf{x}_0,t_0)\Big].$$

- $P(\mathbf{x},t\,|\,\mathbf{x}_0,t_0) \triangleq \text{Prob}\{\mathbf{X}(t)=\mathbf{x},\text{ given }\mathbf{X}(t_0)=\mathbf{x}_0\}$ for $t \ge t_0$.
- Follows from the *probability* statement

$$P(\mathbf{x},t+dt\,|\,\mathbf{x}_0,t_0) = P(\mathbf{x},t\,|\,\mathbf{x}_0,t_0)\times\left[1-\sum_{j=1}^{M}\big(a_j(\mathbf{x})dt\big)\right]$$

$$+\sum_{j=1}^{M}P(\mathbf{x}-\boldsymbol{\nu}_j,t\,|\,\mathbf{x}_0,t_0)\times\big(a_j(\mathbf{x}-\boldsymbol{\nu}_j)dt\big).$$

- But the CME is usually too hard to solve.

---

- Averages: $\quad\big\langle f\big(\mathbf{X}(t)\big)\big\rangle \triangleq \sum_{\mathbf{x}} f(\mathbf{x})P(\mathbf{x},t\,|\,\mathbf{x}_0,t_0)$.

  If we multiply the CME through by $\mathbf{x}$ and then sum over $\mathbf{x}$, we find

$$\frac{d\big\langle \mathbf{X}(t)\big\rangle}{dt} = \sum_{j=1}^{M}\boldsymbol{\nu}_j\big\langle a_j\big(\mathbf{X}(t)\big)\big\rangle.$$

- *If* there were *no fluctuations*,

$$\big\langle a_j\big(\mathbf{X}(t)\big)\big\rangle = a_j\big(\big\langle \mathbf{X}(t)\big\rangle\big) = a_j\big(\mathbf{X}(t)\big),$$

  and the above would reduce to:

$$\frac{d\mathbf{X}(t)}{dt} = \sum_{j=1}^{M}\boldsymbol{\nu}_j a_j\big(\mathbf{X}(t)\big).$$

  - **This is the reaction-rate equation (RRE).**
  - It's usually written in terms of the *concentration* $\mathbf{Z}(t)\triangleq\mathbf{X}(t)/\Omega$.

  ▪ **But as yet, we have *no justification* for ignoring fluctuations.**

➢ **2.** The ***stochastic simulation algorithm*** (SSA):

A procedure for constructing *sample paths* or *realizations* of $\mathbf{X}(t)$.

    *Idea*: Generate *properly distributed random numbers* for
- the time $\tau$ to the *next* reaction,
- the index $j$ of that reaction.

- $p(\tau, j \mid \mathbf{x}, t)\, d\tau \triangleq$ probability, given $\mathbf{X}(t) = \mathbf{x}$, that the *next* reaction will occur in $[t+\tau, t+\tau+d\tau)$, *and* will be $R_j$.

$$= P_0(\tau) \times a_j(\mathbf{x})\, d\tau, \quad P_0(\tau) \triangleq \Pr(no \text{ reactions in time } \tau).$$

$$P_0(\tau + d\tau) = P_0(\tau) \times \left(1 - a_0(\mathbf{x})\, d\tau\right), \quad \text{where} \quad a_0(\mathbf{x}) \triangleq \sum_1^M a_{j'}(\mathbf{x}).$$

Implies $\dfrac{dP_0(\tau)}{d\tau} = -a_0(\mathbf{x}) P_0(\tau).$     Solution: $P_0(\tau) = e^{-a_0(\mathbf{x})\tau}.$

$$\therefore \ p(\tau, j \mid \mathbf{x}, t) = e^{-a_0(\mathbf{x})\tau}\, a_j(\mathbf{x}) = a_0(\mathbf{x}) e^{-a_0(\mathbf{x})\tau} \times \frac{a_j(\mathbf{x})}{a_0(\mathbf{x})}.$$

Thus,
- $\tau$ is an *exponential random variable* with mean $1/a_0(\mathbf{x})$,
- $j$ is an *integer random variable* with probabilities $a_j(\mathbf{x})/a_0(\mathbf{x})$.

---

### The "Direct" Version of the SSA

**1.** In state **x** at time $t$, evaluate $a_1(\mathbf{x}), \ldots, a_M(\mathbf{x})$, and $a_0(\mathbf{x}) \equiv \sum_{j'=1}^M a_{j'}(\mathbf{x})$.

**2.** Draw two unit-interval uniform random numbers $r_1$ and $r_2$, and compute $\tau$ and $j$ according to

- $\tau = \dfrac{1}{a_0(\mathbf{x})} \ln\left(\dfrac{1}{r_1}\right),$

- $j =$ the *smallest integer* satisfying $\sum_{k=1}^{j} a_k(\mathbf{x}) > r_2\, a_0(\mathbf{x}).$

**3.** Replace $t \leftarrow t + \tau$ and $\mathbf{x} \leftarrow \mathbf{x} + \boldsymbol{\nu}_j$.

**4.** Record $(\mathbf{x}, t)$. Return to Step **1**, or else end the simulation.

*A Simple Example*: $S_1 \xrightarrow{c_1} 0$.

$$a_1(x_1) = c_1 x_1, \quad \nu_1 = -1. \text{ Take } X_1(0) = x_1^0.$$

**RRE**: $\dfrac{dX_1(t)}{dt} = -c_1 X_1(t)$. Solution is $X_1(t) = x_1^0 \, e^{-c_1 t}$.
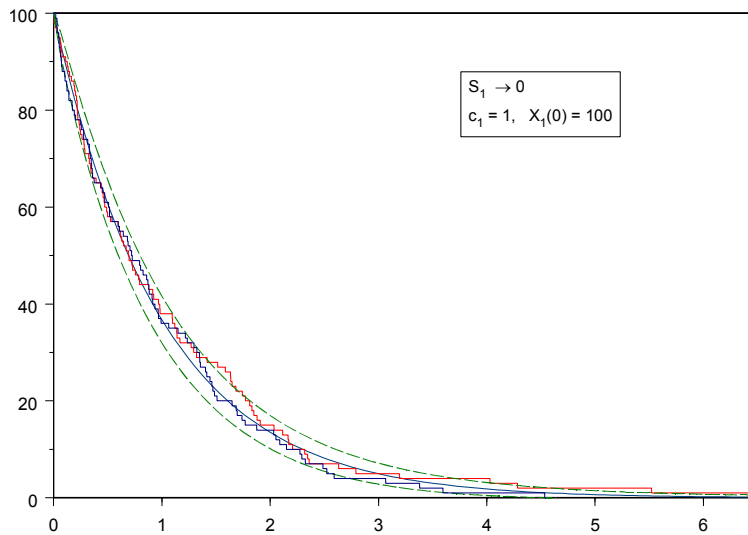
**CME**: $\dfrac{\partial P(x_1, t \mid x_1^0, 0)}{\partial t} = c_1 \big[ (x_1 + 1) P(x_1 + 1, t \mid x_1^0, 0) - x_1 P(x_1, t \mid x_1^0, 0) \big]$.
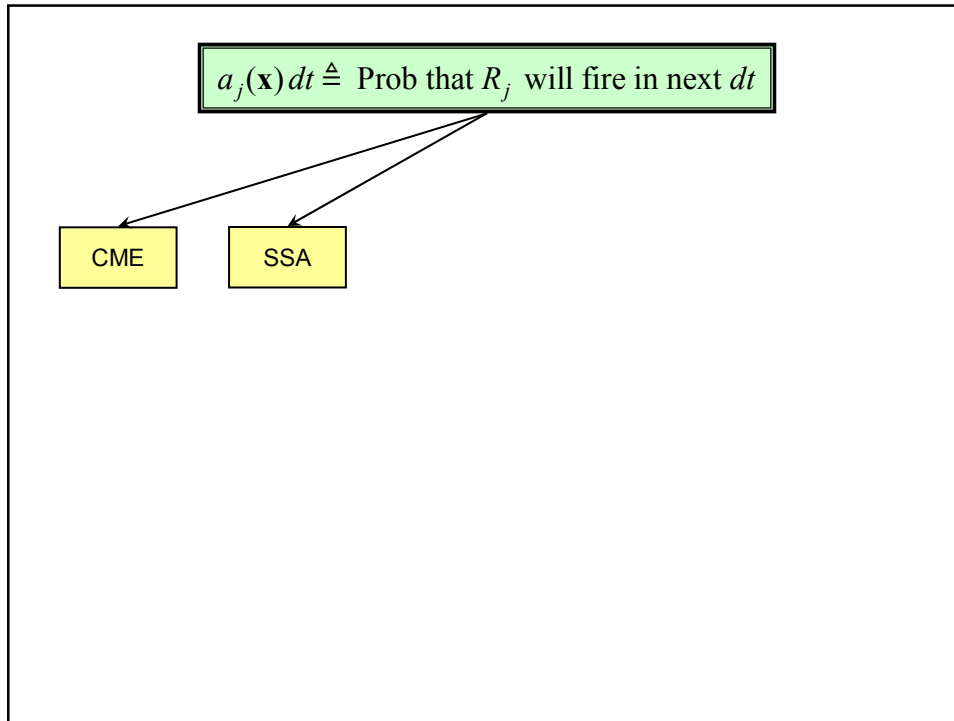
Solution: $P(x_1, t \mid x_1^0, 0) = \dfrac{x_1^0!}{x_1!(x_1^0 - x_1)!} \, e^{-c_1 x_1 t} \left(1 - e^{-c_1 t}\right)^{x_1^0 - x_1} \ (x_1 = 0, 1, \ldots, x_1^0)$

which implies $\langle X_1(t) \rangle = x_1^0 \, e^{-c_1 t}$, $\operatorname{sdev}\{X_1(t)\} = \sqrt{x_1^0 \, e^{-c_1 t} \left(1 - e^{-c_1 t}\right)}$.

**SSA**: Given $X_1(t) = x_1$, generate $\tau = \dfrac{1}{c_1 x_1} \ln\!\left(\dfrac{1}{r}\right)$, then update:

$$t \leftarrow t + \tau, \quad x_1 \leftarrow x_1 - 1.$$



$S_1 \rightarrow 0$
$c_1 = 1, \quad X_1(0) = 100$

$$a_j(\mathbf{x})\,dt \triangleq \text{Prob that } R_j \text{ will fire in next } dt$$

CME

SSA

---

**The SSA . . .**

- Is *exact*.  It does *not* entail approximating "*dt*" by " *Δt* ".
- Is *procedurally simple*, even when the CME is intractable.
- Comes in a variety of implementations …
    - Direct Method (Gillespie, 1976)
    - First Reaction Method (Gillespie, 1976)
    - Next Reaction Method (Gibson & Bruck, 2000)
    - First Family Method (Lok, 2003)
    - Modified Direct Method (Cao, Li & Petzold, 2004)
    - Sorting Direct Method (McCollum, et al. 2006)
- ***Remains too slow for most practical problems***:  Simulating *every* reaction event *one* at a time just takes too much time if any reactants are present in very large numbers.

**We would be willing to sacrifice a little exactness . . .**
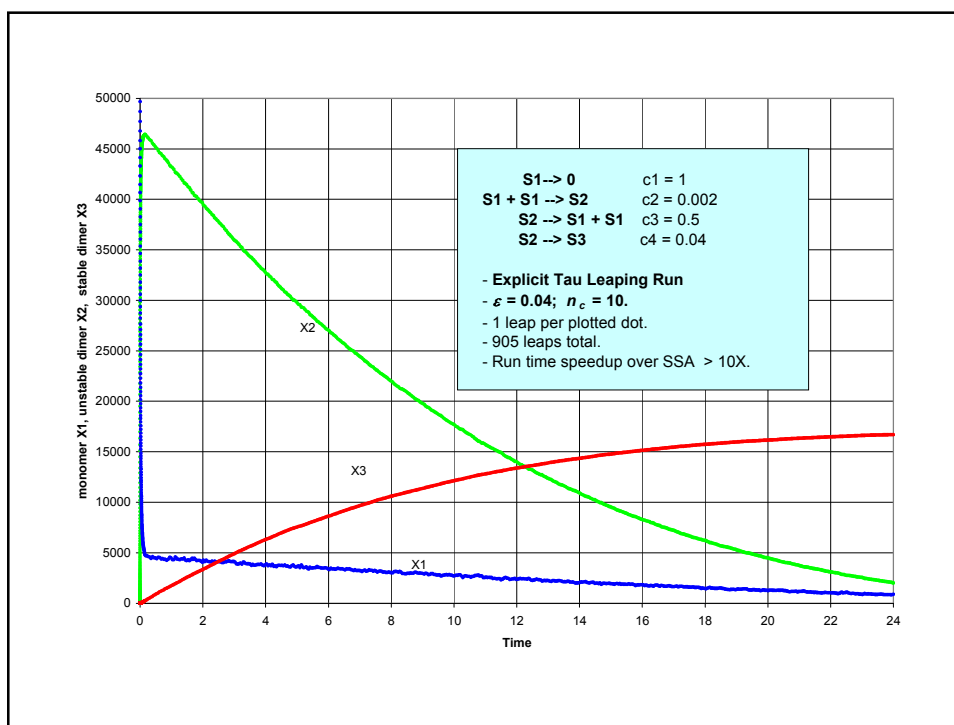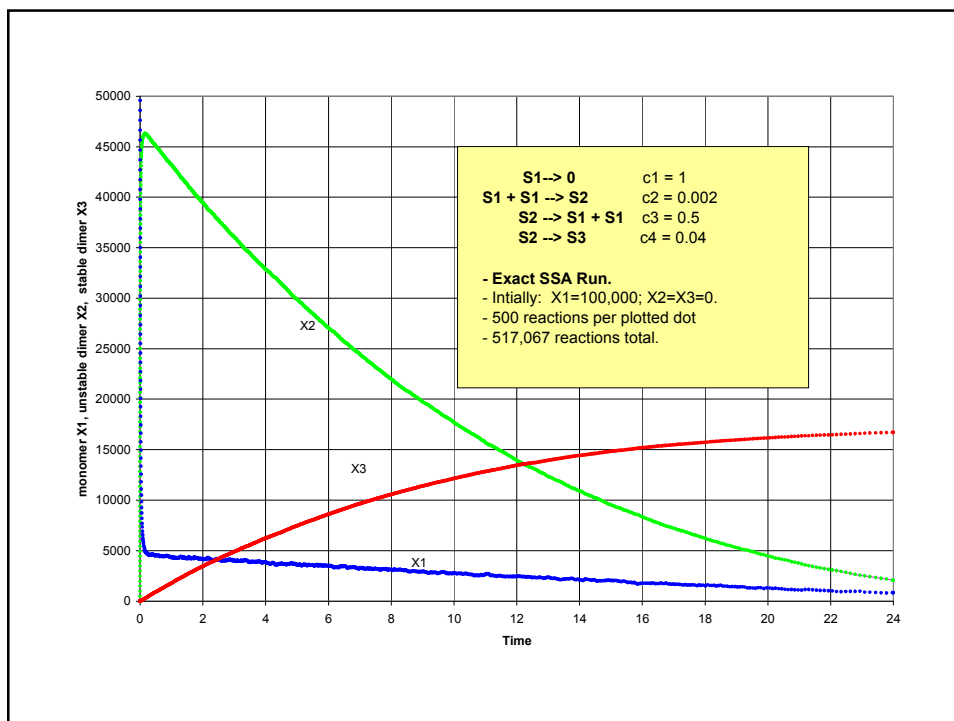**. . . if that would buy us a faster simulation.**

*Tau-Leaping*

- *Approximately* advances the process by a *pre-selected* time $\tau$, which may encompass *more than one* reaction event.

- *Key*: The "Poisson random variable with mean $a\tau$" can be defined:
  $\mathcal{P}(a\tau) \equiv$ the **number of events** that will occur in a time $\tau$,
  when the probability of an event in any $dt$ is $adt$,
  provided $a$ is a positive **constant**.

- With $\mathbf{X}(t) = \mathbf{x}$, let us choose $\tau$ *small enough* to satisfy the
  **Leap Condition**: Each $a_j(\mathbf{x}) \approx$ **constant** in $[t, t+\tau]$.

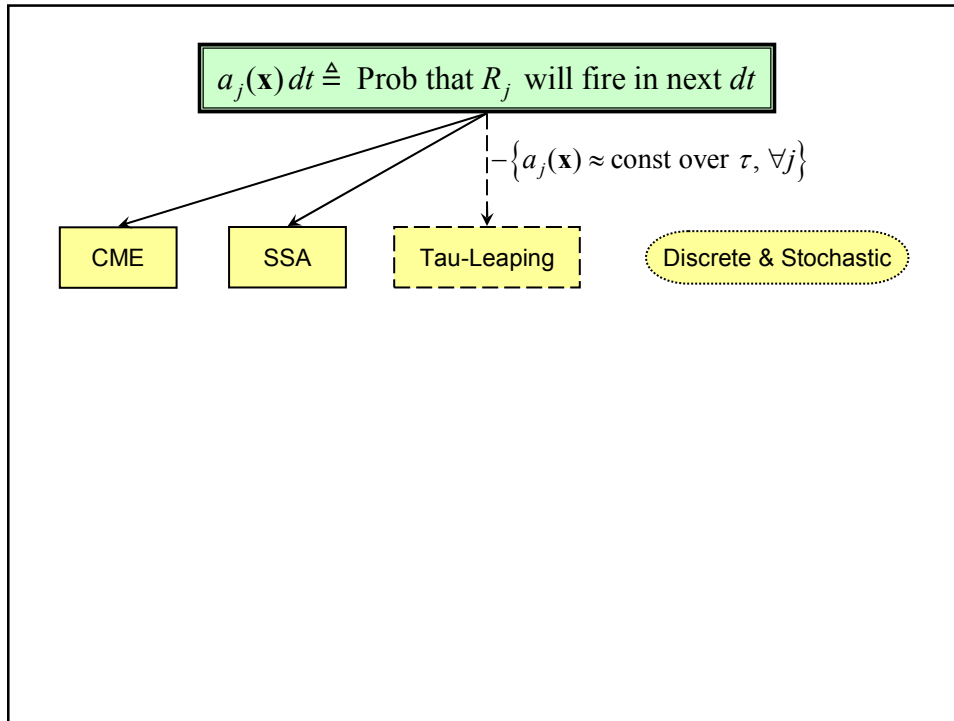- Then the number of $R_j$ firings in $[t, t+\tau]$ will be $\approx \mathcal{P}(a_j(\mathbf{x})\tau)$. So ...

$$\mathbf{X}(t+\tau) \doteq \mathbf{x} + \sum_{j=1}^{M} \mathcal{P}_j\left(a_j(\mathbf{x})\tau\right) \boldsymbol{\nu}_j$$

---

$$\mathbf{X}(t+\tau) \doteq \mathbf{x} + \sum_{j=1}^{M} \mathcal{P}_j\left(a_j(\mathbf{x})\tau\right) \boldsymbol{\nu}_j$$

*- Practical Considerations for Implementing Tau-Leaping -*

- **Finding the largest $\tau$ that satisfies the Leap Condition.**
  - Accomplished via an accuracy control parameter $\varepsilon$.
  - We estimate the largest $\tau$ for which $\left|\Delta_\tau a_j / a_j\right| \le \varepsilon, \forall j$.

- **Avoiding negative populations, and segueing to the SSA**.
  - Accomplished via a second control parameter $n_c$.
  - We call any reaction that is within $n_c$ firings of exhausting any reactant a *critical* reaction. Then we tau-leap *no farther* than the *next* firing of a critical reaction.
  - Becomes the SSA when **all** reactions are classified critical.

Chart 1 text content:

monomer X1, unstable dimer X2, stable dimer X3

| S1 --> 0 | c1 = 1 |
| S1 + S1 --> S2 | c2 = 0.002 |
| S2 --> S1 + S1 | c3 = 0.5 |
| S2 --> S3 | c4 = 0.04 |

- **Exact SSA Run.**
- Intially:  X1=100,000; X2=X3=0.
- 500 reactions per plotted dot
- 517,067 reactions total.

X2, X3, X1

Time



Chart 2 text content:

monomer X1, unstable dimer X2, stable dimer X3

| S1 --> 0 | c1 = 1 |
| S1 + S1 --> S2 | c2 = 0.002 |
| S2 --> S1 + S1 | c3 = 0.5 |
| S2 --> S3 | c4 = 0.04 |

- **Explicit Tau Leaping Run**
- $\varepsilon = 0.04$;  $n_c = 10$.
- 1 leap per plotted dot.
- 905 leaps total.
- Run time speedup over SSA  > 10X.

X2, X3, X1

Time

$$a_j(\mathbf{x})\,dt \triangleq \text{Prob that } R_j \text{ will fire in next } dt$$

$$\vdash\{a_j(\mathbf{x}) \approx \text{const over } \tau,\ \forall j\}$$

| CME | SSA | Tau-Leaping | Discrete & Stochastic |

---

**Speeding up Tau-Leaping: The Langevin Equation**

- Two math facts:
    - If $m \gg 1$, then $\mathcal{P}(m) \approx \mathcal{N}(m,m)$.
    - $\mathcal{N}(m,\sigma^2) = m + \sigma\mathcal{N}(0,1)$.

- So, with $\mathbf{X}(t) = \mathbf{x}$, suppose we can choose $\tau$ *small enough* to satisfy the Leap Condition, *yet also large enough that* $a_j(\mathbf{x})\tau \gg 1,\ \forall j$.

Then . . . 
$$\mathbf{X}(t+\tau) \doteq \mathbf{x} + \sum_{j=1}^{M} \mathcal{P}_j\big(a_j(\mathbf{x})\tau\big)\boldsymbol{\nu}_j$$

$$\doteq \mathbf{x} + \sum_{j=1}^{M} \mathcal{N}_j\big(a_j(\mathbf{x})\tau, a_j(\mathbf{x})\tau\big)\boldsymbol{\nu}_j$$

$$\doteq \mathbf{x} + \sum_{j=1}^{M} \Big[a_j(\mathbf{x})\tau + \sqrt{a_j(\mathbf{x})\tau}\,\mathcal{N}_j(0,1)\Big]\boldsymbol{\nu}_j$$

❖ $\quad \mathbf{X}(t+\tau) \doteq \mathbf{x} + \sum_{j=1}^{M}\boldsymbol{\nu}_j a_j(\mathbf{x})\tau + \sum_{j=1}^{M}\boldsymbol{\nu}_j\sqrt{a_j(\mathbf{x})}\,\mathcal{N}_j(0,1)\sqrt{\tau}$ .

$$\mathbf{X}(t+\tau) \doteq \mathbf{x} + \sum_{j=1}^{M} \boldsymbol{\nu}_j \, a_j(\mathbf{x}) \tau + \sum_{j=1}^{M} \boldsymbol{\nu}_j \sqrt{a_j(\mathbf{x})} \, \mathcal{N}_j(0,1) \sqrt{\tau}$$

- This is the ***Langevin leaping formula***.
- It's faster than the ordinary tau-leaping formula, because
  - $a_j(\mathbf{x})\tau \gg 1$ means *lots* of reaction events get leapt over in $\tau$ ;
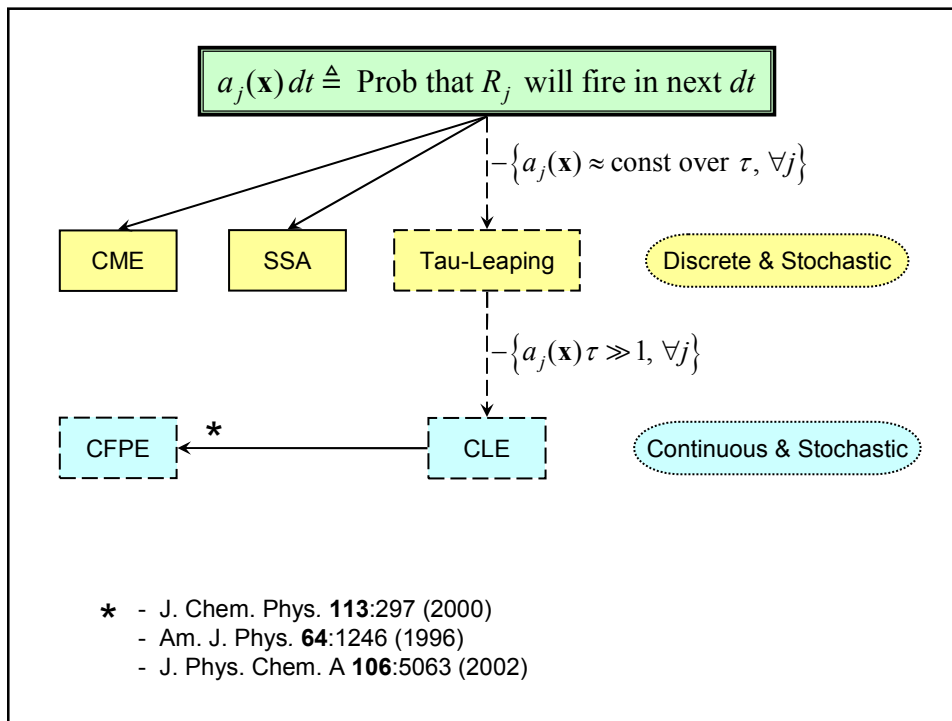  - *normal* random numbers can be generated faster than Poissons.
- It directly implies, and is entirely equivalent to, a SDE called
  the ***chemical Langevin equation*** (CLE):

$$\frac{d\mathbf{X}(t)}{dt} \doteq \sum_{j=1}^{M} \boldsymbol{\nu}_j \, a_j(\mathbf{X}(t)) + \sum_{j=1}^{M} \boldsymbol{\nu}_j \sqrt{a_j(\mathbf{X}(t))} \, \Gamma_j(t) \ .$$

  - *Gaussian white noise*: $\Gamma(t) \triangleq \lim_{dt \to 0^+} \dfrac{\mathcal{N}(0,1)}{\sqrt{dt}} \equiv \lim_{dt \to 0^+} \mathcal{N}\left(0, \dfrac{1}{dt}\right).$
  - Satisfies $\left\langle \Gamma_j(t)\, \Gamma_{j'}(t') \right\rangle = \delta_{jj'}\, \delta(t-t').$
- Our *discrete stochastic* process $\mathbf{X}(t)$ has now been *approximated* as a
  *continuous stochastic* process.



$a_j(\mathbf{x})\,dt \triangleq$ Prob that $R_j$ will fire in next $dt$

$\vdash \left\{ a_j(\mathbf{x}) \approx \text{const over } \tau,\ \forall j \right\}$

CME   SSA   Tau-Leaping   Discrete & Stochastic

$\vdash \left\{ a_j(\mathbf{x}) \tau \gg 1,\ \forall j \right\}$

CFPE  ←*  CLE   Continuous & Stochastic

* - J. Chem. Phys. **113**:297 (2000)
  - Am. J. Phys. **64**:1246 (1996)
  - J. Phys. Chem. A **106**:5063 (2002)

<div style="border:1px solid black; padding:10px;">

**The Thermodynamic Limit**

***Def:*** All $X_i \to \infty$, and $\Omega \to \infty$, with $X_i/\Omega$ constants.

- $\left.\begin{array}{l} a_j = c_j x_1 \sim x_1 \\ a_j = c_j x_1 x_2 \sim \Omega^{-1} x_1 x_2 \sim x_2 \end{array}\right\} \Rightarrow \left\{\begin{array}{l} \text{In the thermodynamic limit,} \\ \textbf{all } a_j\text{'s grow like (system size).} \end{array}\right.$

- So in the thermodynamic limit, we see that in the CLE

$$\frac{d\mathbf{X}(t)}{dt} \doteq \sum_{j=1}^{M} \boldsymbol{\nu}_j\, a_j\big(\mathbf{X}(t)\big) + \sum_{j=1}^{M} \boldsymbol{\nu}_j \sqrt{a_j\big(\mathbf{X}(t)\big)}\, \Gamma_j(t) \ ,$$
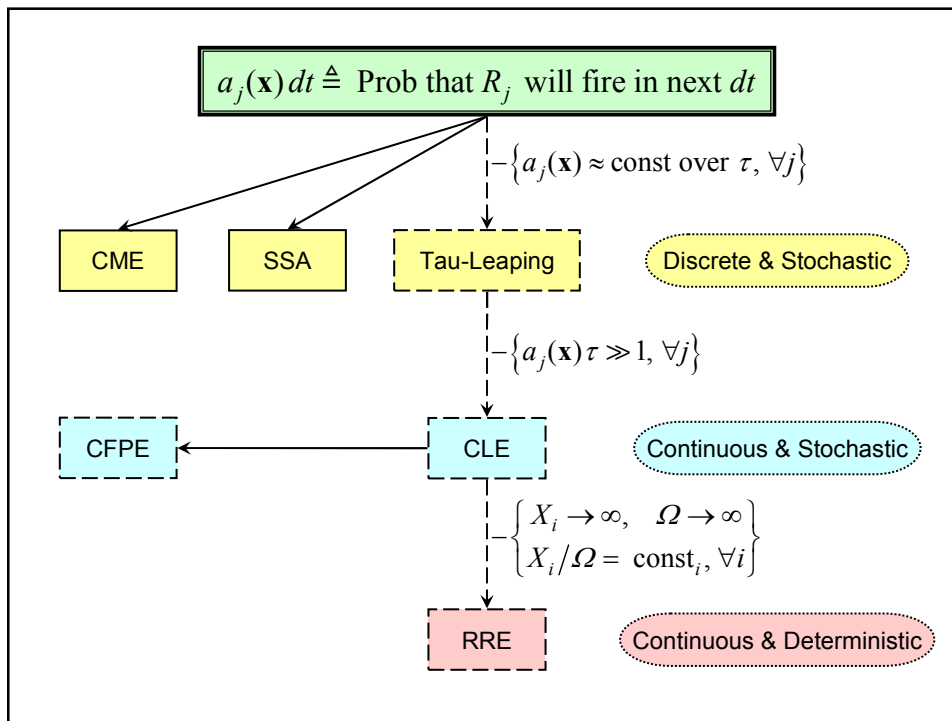
  - the *deterministic* term grows like (system size),
  - the *stochastic* term grows like $(\text{system size})^{1/2}$.

- $\Rightarrow$ Rule of Thumb: ***Relative fluctuations die off as (system size)$^{-1/2}$.***

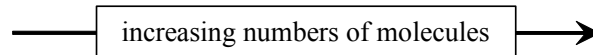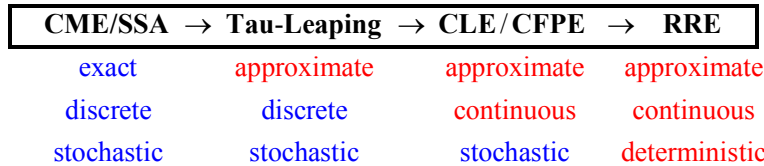- **At the thermodynamic limit** the stochastic term disappears, leaving

$$\frac{d\mathbf{X}(t)}{dt} \doteq \sum_{j=1}^{M} \boldsymbol{\nu}_j\, a_j\big(\mathbf{X}(t)\big) \ \ \dots \text{ the } \textbf{RRE \dots \textit{derived!}}$$

  $\mathbf{X}(t)$ has now become a ***continuous deterministic*** process.

</div>

***A Spectrum of Descriptive Mathematical Modes***
(for **well-stirred** systems)

$$a_j(\mathbf{x}) \approx \text{const} \qquad a_j(\mathbf{x})\,\tau \gg 1 \qquad X_i \to \infty,\; \Omega \to \infty$$

$$\underbrace{\text{in } [t, t+\tau],\, \forall j} \qquad \underbrace{\forall j} \qquad \underbrace{X_i / \Omega \to \text{const}_i}$$

| | | |
|:---:|:---:|:---:|
| | | |

$$\boxed{\text{CME/SSA} \;\to\; \text{Tau-Leaping} \;\to\; \text{CLE}/\text{CFPE} \;\to\; \text{RRE}}$$

| exact | approximate | approximate | approximate |
|:---:|:---:|:---:|:---:|
| discrete | discrete | continuous | continuous |
| stochastic | stochastic | stochastic | deterministic |

→ increasing numbers of molecules →

---

### *Another Multi-Scale Problem*

- Some reactions/species may be *very fast*, others *very slow*.
- "Fast" and "slow" are **interconnected** – not easy to separate.
- Often manifests as *dynamical stiffness*, a known ODE problem.
- SSA still works, and is exact. But it's agonizingly slow.
- Tau-leaping remains accurate, but the Leap Condition restricts $\tau$ to the shortest (fastest) time scale of the system. So even it's too slow.
- ***One approach:*** *Implicit Tau-Leaping*
  **-** A stochastic adaptation of the *implicit Euler method* for ODEs.
- ***Another approach:*** The *Slow-Scale Stochastic Simulation Algorithm*
  **-** Skips over the fast reactions and *simulates only the slow reactions*, using specially modified propensity functions. An adaptation of the "rapid equilibrium" / "quasi steady-state" methods for RREs.