

Supplementary Material

Details of the algorithm for the GPU architecture

Algorithm

The algorithm for the GPU architecture was adapted from the Gillespie implementation of Mauch et. al. which is, to our knowledge, the fastest single core implementation available.¹ The major differences of the algorithm are the utilization of the ziggurat method to generate the exponential deviates combined with a 2D search on the propensities and continuous updating.

The ziggurat implementation was taken from the original publication with the minor difference that we did not use the proposed SHR3 generator but rather a concatenated multiply-with-carry generator (MWC) with random seeding which passes the DIEHARD tests for randomness.^{2,3} This has virtually no impact on the speed of the random number generation. The 2D search is implemented as in the Mauch version. The propensities are stored in a matrix with row and column size equal to the smallest integer larger than the square root of the reaction count. For each row we also store the row sum of propensities. The uniform variate used for the search is also generated from the MWC. The 2D search first identifies the row of the corresponding propensity row sum and then looks within the row for the corresponding propensity. After identifying the next reaction updates are only executed for propensities influenced by the executed reaction. This is achieved by automatic computation of a dependency list prior to the simulation from the stoichiometry of the model. The propensity function itself is compiled together with the simulator. For a schematic view of the algorithm also see Fig. S1B.

Parallelization and memory management

In order to efficiently parallelize the algorithm we adapted the memory accesses in the model to comply with NVIDIA hardware. All matrix data structures used during the simulation are linearized first to meet alignment constraints and access patterns have been optimized. The data structures used during the simulation (dependencies, stoichiometry, propensities and state) are kept within shared memory, which is on chip and can be accessed two orders of magnitude faster than global memory. Every single thread will now simulate a single trajectory. Additionally, every single thread will only access a single memory bank in shared memory, thus minimizing the number of bank conflicts. In particular, there will be no bank conflicts at all if the number of threads is equal or less to the number of banks (32 for devices of compute capability 2.x). Recording in the state is done in

global memory with a maximum amount of coalescing and caching. As such, recording of the state can be bundled into one write operation for several threads (also see Fig. S1A). The tables for the exponential random number generation are kept within constant memory, and are, thus, always cached.

As random number generation using the described methods is rapid on the GPU and since the propensity calculation uses only simple float operations simulation speed is mostly bandwidth-limited. Thus, we chose the thread number as a compromise between a relatively low thread number (to avoid bank conflicts) and maximum occupancy.

Benchmark

The algorithm was benchmarked on a NVIDIA Geforce 480GTX. The chosen thread number per block was 192 which resulted in occupancy of 75%. The benchmark model was the Decaying Dimerization model which comprises the following reactions r_i :



with corresponding parameters $k_1 = 1.0$, $k_2 = 0.002$, $k_3 = 0.5$ and $k_4 = 0.04$.

The average simulation time per reaction for our method and two other implementations is shown in Fig. S2. For implementations other than our own the average time per reaction has been calculated from the mean and the extreme values for the number of reactions per trajectory. The mean, minimum and maximum number of reactions were calculated from a set of 20 million trajectories. The CPU times have been generated with Cain on Dual Core Intel processor (blue line). This corresponds perfectly to the minimal single core CPU simulation time of about 105 ns for a single reaction. The green lines denote a GPU implementation of Li et. al.⁴ The observed jump of simulation is likely to be caused by a switch in the thread number which increases occupancy of the GPU. However, since their simulation times come from NVIDIA GeForce 8800GTX which has a different hardware setup as well as a different compute capability we cannot say what exactly causes the speed up of our implementation compared to theirs. The red line denotes our own implementation which always shows the fastest simulation times which peak at a maximum speed of 0.53 ns per reaction (more than 1.8 billion reactions per second). A simulation speed of around 0.6 ns per reaction is achieved from 10,000 trajectories on and thus makes it applicable for a large variety of tasks. Additionally, recording of several time steps only mildly affects simulation times.

References

- 1 S. Mauch and M. Stalzer, *IEEE/ACM Trans. Comput. Biol. Bioinform.*, 2011, **8**, 27–35.
- 2 G. Marsaglia, *J. Statist. Soft.*, 2003, **8**, 1–6.
- 3 G. Marsaglia and W. W. Tsang, *J. Statist. Soft.*, 2000, **5**, 1–7.
- 4 H. Li and L. Petzold, *IJHPCA*, 2009, **24**, 107–116.

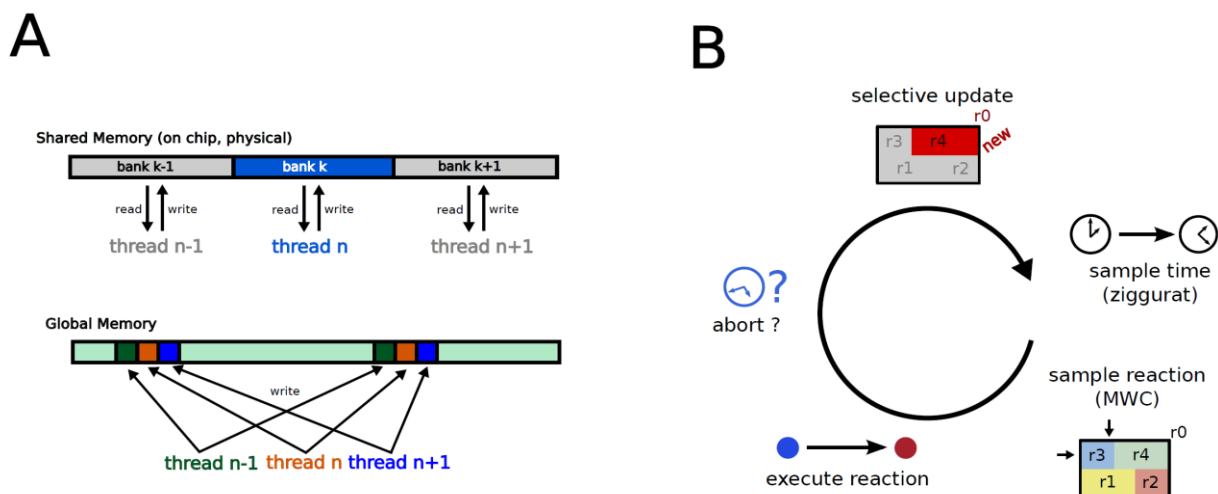


Fig. S1 Schematic view of the memory management (A) and the implemented algorithm (B).

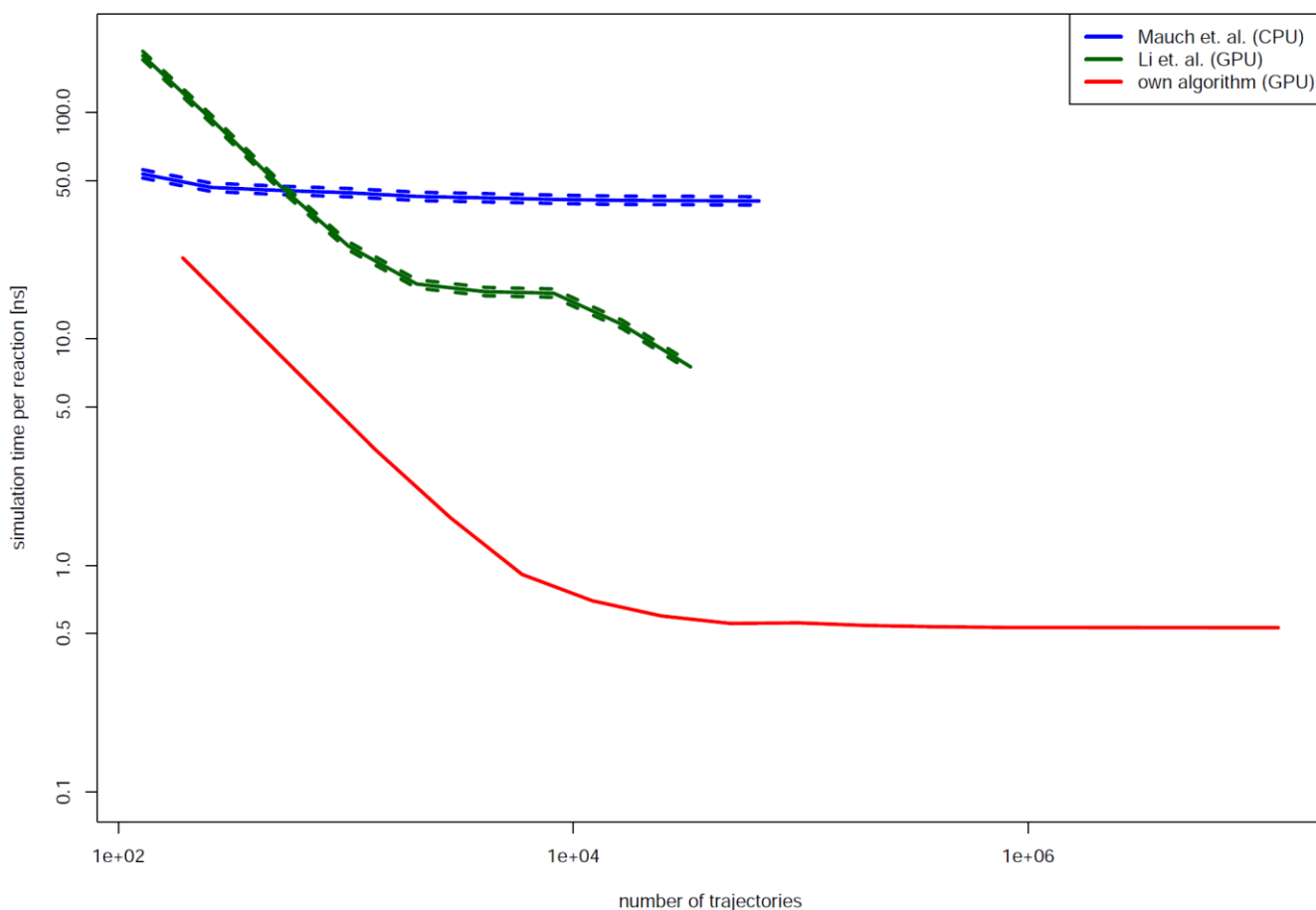


Fig. S2 Comparison of simulation times per reaction for the Decaying Dimerization model. The CPU used was a Intel Core 2 Duo 3.16 GHz. The GPUs used were a NVIDIA GeForce 8800 GTX (Li. et. al.) and a NVIDIA GeForce 480 GTX (this publication). The dashed lines correspond to minimum and maximum approximations.

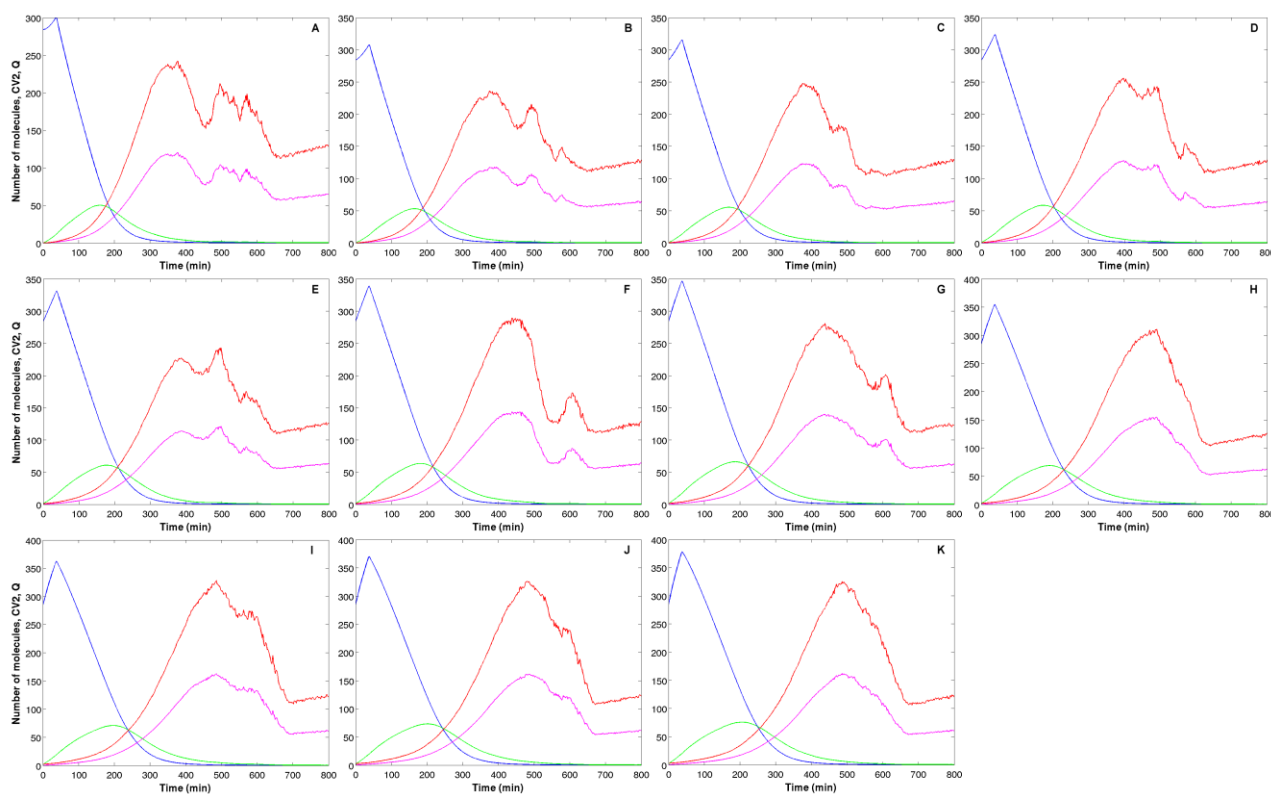


Fig. S3 Simulated dynamics of mean (blue) and standard deviation (green) of Sic1, of CV2 (pink) and Q (red) for a range from 0 (A) to 10 (K) initial *SIC1* mRNA molecules at a ratio of k_1/k_2 equal to 3. The curves of CV2 and Q were multiplied by 50 for purpose of visualization.

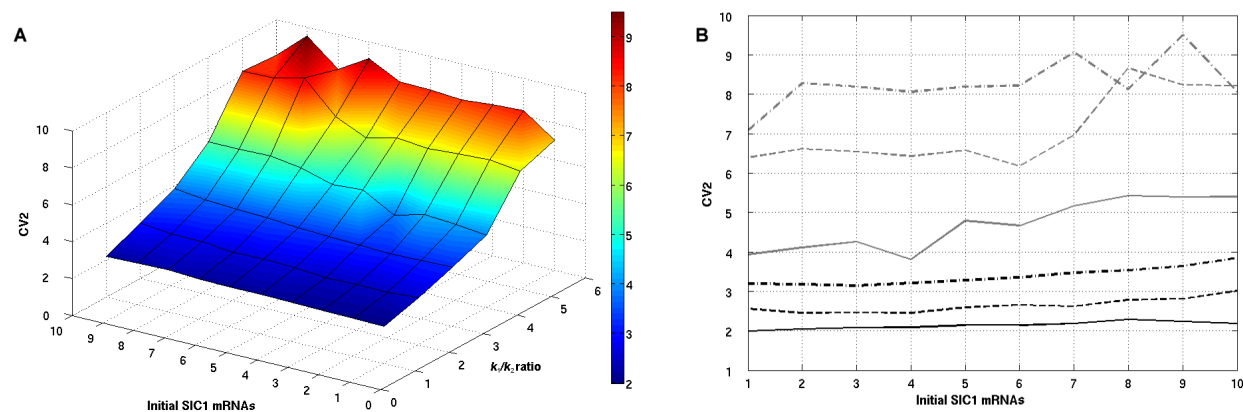


Fig. S4 Relation between CV2, the initial *SIC1* mRNA molecule number and different k_1/k_2 ratios of *SIC1* mRNA shown with a three-dimensional representation (A) and a bi-dimensional representation (B). In panel B, dash-dotted, dashed and solid lines represent the ratios k_1/k_2 from 1 to 6.

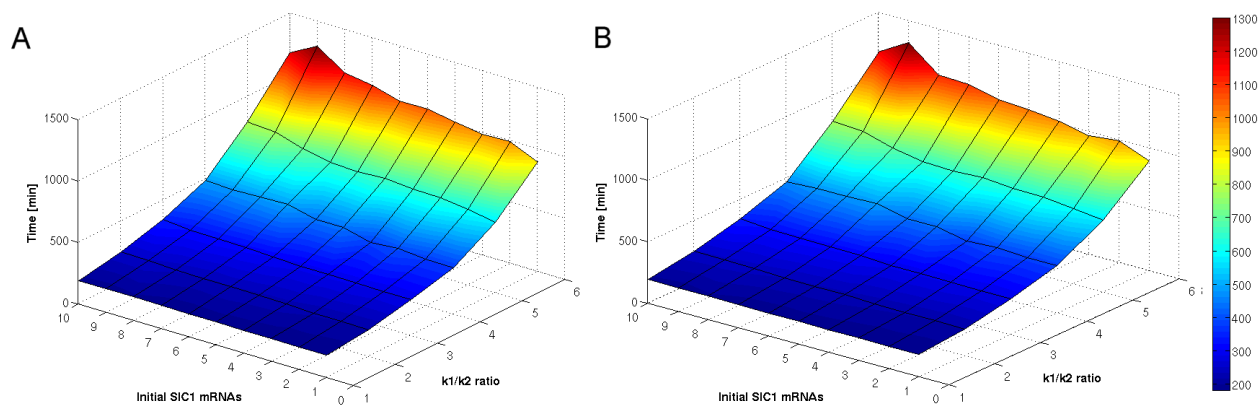


Fig. S5 Relation between the Time of peaks for Q (A) and CV2 (B), the initial *SIC1* mRNA molecule number and different k_1/k_2 ratios of *SIC1* mRNA shown with a three-dimensional representation.

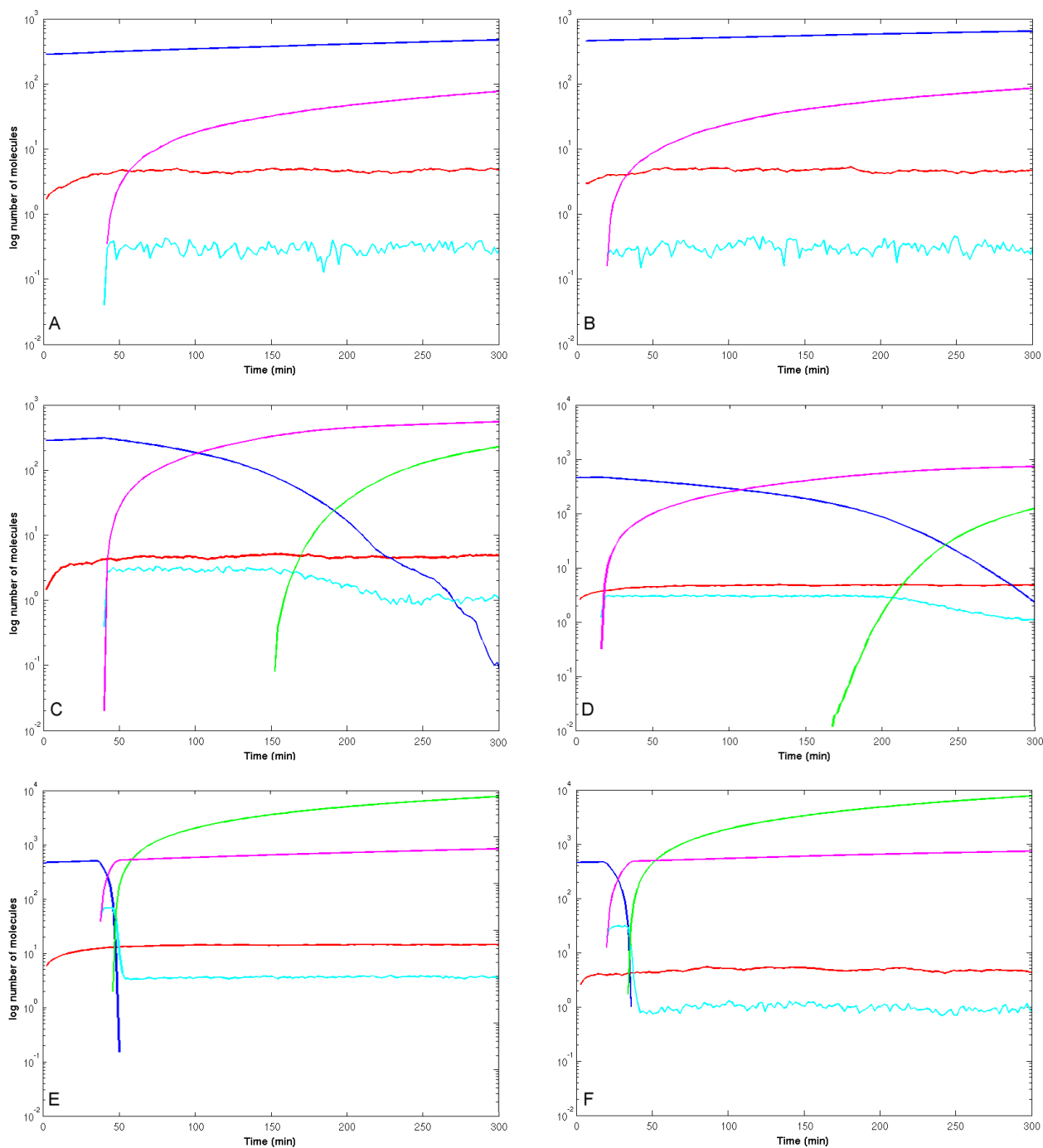


Fig. S6 Semi-log plots of data shown in Fig. 6. Simulated dynamics of a daughter cell (left column) and a mother cell (right column) for different Clb5,6 production rate constants: k_{6a} (A, B), k_{6b} (C, D) and k_{6c} (E, F). Sic1 (blue), cytoplasmic Clb5,6 (green), Sic1-Clb5 complex (light blue), nuclear Clb5,6 (pink) and *SIC1* mRNA (red) are shown.

Table S1 Chemical notation of the reactions of the G1/S transition stochastic model

Reaction	ID	Chemical Notation
re_1	Prod_mRNA	\rightarrow mRNA_Sic1
re_2	Deg_mRNA	mRNA_Sic1 \rightarrow
re_3	Prod_Sic1	mRNA_Sic1 \rightarrow mRNA_Sic1 + Sic1
re_4	Formation_of_Complex	Sic1 + Clb5 \rightarrow Sic1Clb5
re_5	Decay_of_Complex	Sic1Clb5 \rightarrow Sic1 + Clb5_active
re_6 (a, b, c)	Prod_Clb5	\rightarrow Clb5
re_7	Deg_Sic1	Sic1 \rightarrow

Table S2 Initial conditions for a single daughter cell

Species ID	Amount	
mRNA_SIC1	varied between 0 and 12	
Sic1	738	
Sic1Clb5	0	
Clb5_active	0	
Clb5	0	
Reaction (constant)	Rate Constant Value (Deterministic)	Rate Constant Value (Stochastic)
re_1 (k_1)	0.1 min^{-1}	0.1 min^{-1}
re_2 (k_2)	varied between: 0.1 min^{-1} , 0.05 min^{-1} , 0.03333 min^{-1} , 0.025 min^{-1} , 0.02 min^{-1} , 0.01667 min^{-1} , 0.0143 min^{-1} and 0.0111 min^{-1}	varied between: 0.1 min^{-1} , 0.05 min^{-1} , 0.03333 min^{-1} , 0.025 min^{-1} , 0.02 min^{-1} , 0.01667 min^{-1} , 0.0143 min^{-1} and 0.0111 min^{-1}
re_3 (k_3)	0.32 min^{-1}	0.32 min^{-1}
re_4 (k_4)	$84.6 \mu\text{M}^{-1} * \text{min}^{-1}$	0.0056 min^{-1}
re_5 (k_5)	1 min^{-1}	1 min^{-1}
re_6 (k_{6b})	0.3 min^{-1}	3.0 min^{-1} (assuming 10 <i>CLB5</i> mRNA molecules)

Table S3 Initial conditions for mother and daughter cells

	Daughter	Mother
t_{G1} (time at the initiation of Clb5,6)	37 min	15.6 min
Cell volume (V_0)	25 fl	40 fl
Initial molecule numbers		
mRNA_SIC1	1	2
Sic1	284	454
Clb5	0	0
Sic1Clb5	0	0
Clb5_active	0	0
Recalculation of rate constants (min^{-1})		
Prod_mRNA (k_1)	0.1	0.1
Deg_mRNA (k_2)	0.03333	0.03333
Prod_Sic1 (k_3)	0.32	0.32
Formation_of_Complex (k_4)	0.0035	0.0056
Decay_of_Complex (k_5)	1	1
Prod_Clb5 (k_{6a} for ~ 1 <i>CLB5</i> mRNAs)	0.3	0.3
Prod_Clb5 (k_{6b} for ~ 10 <i>CLB5</i> mRNAs)	3	3
Prod_Clb5 (k_{6c} for ~ 100 <i>CLB5</i> mRNAs)	30	30
Calculations for reaction re_4		
V_0 / t_{G1}	0.68 fl min^{-1}	2.6 fl min^{-1}