# PERFORMANCE MODELING AND CONTROL OF STRUCTURALLY CONTROLLED RESOURCE ALLOCATION SYSTEMS

## Jin Young Choi and Spyros Reveliotis

*School of Industrial & Systems Engineering*
*Georgia Institute of Technology*
*765 Ferst Drive, Atlanta, GA 30332*

Abstract: While the structural and performance-oriented control problems related to the real-time management of the resource allocation systems (RAS) underlying the operation of many contemporary technological applications have, both, been investigated extensively in the past, their integration in a seamless modeling and analysis framework, and the effects arising from their interaction, remain yet to be addressed. This paper undertakes these issues by investigating the class of Generalized Stochastic Petri nets(GSPN's) as a convenient and powerful analytical framework for the integrated modeling of the logical and the time-based RAS dynamics. More specifically, it is shown that GSPN's provide an exact formulation for the problem of performance control of structurally controlled RAS, in the case of systems with exponentially distributed event firing times, and effective approximations for the more general case of systems with non-exponential event timings. On the theoretical side, the aforementioned formulation can be utilized to (re-)establish various properties of the optimal scheduling policy in the considered operational context, including the existence of an optimal deterministic stationary scheduling policy. Finally, the proposed framework can support the thorough characterization of the structure of the optimal scheduling policy for small-sized RAS's, under various parameterizations, allowing, thus, for a more systematic study and a more profound understanding of the timed dynamics taking place in these environments. *Copyright©2002IFAC*

Keywords: Resource Allocation Systems, structural control, Performance control, scheduling, Generalized Stochastic Petri-nets

## 1. INTRODUCTION

With the migration of modern technological applications to highly automated modes of operation, the effective and efficient deployment, reconfiguration and control of the resource allocation underlying the emerging modes of these environments is an issue of ever-increasing importance. Yet, currently we are lacking an adequate methodology for the effective real-time management of these flexibly automated *resource allocation systems (RAS)*, partly due to the fact that past research on the performance modeling and control of these environments has adopted a high-level perspective of their dynamics, ignoring the lower-level operational details. Characteristically, it is interesting to notice that the *hierarchical decomposition* framework (Gershwin, 1989), the most widely adopted "analytical" framework for planning and control in production environments, discerns strategic, tactical, and operational de-

cisions, all of which address performance objectives, while it presumes the logically consistent and robust system behavior. However, in an extensively automated environment, the establishment of logically correct and robust behavior to the various operational contingencies, is definitely a responsibility of the underlying control logic. This gives rise to a new set of control problems, referred to as the RAS *Structural Control (SC)*, which has been investigated extensively in the last decade, but the integration of the developed set of results with the complementary function of performance-oriented control remains still to be addressed.

The work presented in this paper undertakes the development of an analytical framework for performance modeling and control of *structurally controlled RAS*, that will support the seamless integration of logical and performance-oriented control problems. More specifically, it is shown that the class of the *Generalized Stochastic Petri nets (GSPN's)* offers a convenient and powerful analytical framework for modeling the considered problem, providing exact formulations in the case of systems with exponentially distributed event firing times, and effective approximations for the more general case of systems with non-exponential event timings. Furthermore, on the theoretical side, by providing a detailed, closed-form representation of the optimal scheduling problem for structurally controlled RAS, the aforementioned GSPN modeling and analysis framework can be used in order to (re-)establish some properties regarding the structure of the optimal scheduling policies; in particular, it can be shown that, for a large part of the considered class of problems, there will always exist an optimal *deterministic stationary* scheduling policy. Finally, for the case of small-sized systems, the proposed framework supports the thorough characterization of the structure of the optimal policies under various system parameterizations, allowing, thus, for a more systematic study and a more profound understanding of the (timed) dynamics taking place in these environments. The rest of the paper is organized as follows: Section 2 introduces the proposed GSPN-based analytical framework for performance modeling and analysis of structurally controlled RAS. Section 3 discusses the existence of a deterministic stationary optimal performance control policy for the considered problem. Section 4 elucidates the concepts and the analytical power of the proposed modeling framework, by undertaking the detailed analysis of a small re-entrant line with finite buffering capacity. Finally, Section 5 concludes the paper, and suggests directions for future work. In the following, it is assumed that the reader is familiar with the GSPN model and the basic RAS theory; an excellent introduction of the former can be found in (A. Marsan

*et al.*, 1986), while an overview of the latter is provided at (Reveliotis, 2001).

## 2. GSPN-BASED PERFORMANCE MODELING AND ANALYSIS

The key idea underlying the proposed approach for modeling the time-based dynamics of the structurally controlled RAS through the class of GSPN's, is to refine the resource-process nets that typically model the qualitative/logical behavior of the considered RAS (Banaszak and Krogh, 1990), by (i) introducing a more detailed modeling of the processing, staging and transport phases experienced by the different job instances during their advancement through their various processing stages, and (ii) explicitly associating a timing distribution with the various transitions corresponding to the job active processing or transfers. Furthermore, in order to facilitate the GSPN-based modeling and analysis, it is assumed that all transition times are exponentially distributed. This last assumption can be relaxed, whenever it is deemed as too unrealistic, by substituting each timed transition in the resulting GSPN model with a GSPN subnet, modeling a phase-type distribution that approximates the original/empirical distribution of the underlying event timings; we refer the reader to (Papadopoulos *et al.*, 1993) for a detailed treatment of phase-type distributions and the relevant approximation theory. The details underlying the GSPN-based modeling of structurally controlled RAS are further elucidated through the example of Section 4.

According to the general GSPN theory (A. Marsan *et al.*, 1986), the marking process of a GSPN net, $\mathcal{N}$, is a semi-Markov process with a discrete state space, $\mathcal{S}$, given by the net reachability space $R(N, M_0)$. $\mathcal{S}$ is partitioned to *vanishing* states / markings, $\mathcal{V}$, which enable at least one immediate transition of $\mathcal{N}$, and therefore, they have zero sojourn time, and *tangible* markings, $\mathcal{T}$, which enable only timed transitions, and therefore, they present positive sojourn times. Furthermore, the untimed system dynamics, defined by its transitional patterns among the various states of its reachable state space, are characterized by the, so called, *Embedded Markov Chain (EMC)*, whose branching probabilities, $Q = [q_{kl}]$ are determined by externally specified *(dynamic) random switches*, in case of vanishing markings, and the enabled event *exponential race*, in case of tangible markings. If this EMC is finite-state, homogeneous, and irreducible, it possesses a steady-state distribution.

In the case of GSPNs modeling RAS behavior, the underlying EMC is finite-state and homogeneous, while the imposition of a liveness-enforcing SCP

on the RAS behavior establishes also its irreducibility. Furthermore, the set of vanishing markings, $\mathcal{V}$, represents the decision-making points in the RAS operation, while the associated set of *dynamic* random switches, $\Xi$, implements the logic of the imposed scheduling policy. Hence, the set of random switching probabilities, $\xi_l$, defines the decision variables of the underlying RAS scheduling problem, that must be priced in a way that optimizes the performance objective under consideration. Letting $Q(\xi)$ denote the transition probability matrix (TPM) of the system EMC, the problem of optimizing a performance objective for the structurally controlled RAS can be formally expressed by the following Mathematical Programming (MP) formulation:

$$\max_{\xi} \mathcal{O}(\xi) \qquad (1)$$

s.t.

$$\mathbf{y} = \mathbf{y}Q(\xi) \ ; \ \sum_{s_k \in \mathcal{S}} y_k = 1 \qquad (2)$$

$$\pi_k = \begin{cases} 0, & s_k \in \mathcal{V} \\ y_k E[s_k] / \sum_{s_l \in \mathcal{T}} y_l E[s_l] & s_k \in \mathcal{T} \end{cases} \qquad (3)$$

$$\forall s_k \in \mathcal{T}, E[s_k] = 1 \ / \sum_{T_j \text{ enabled in } s_k} r_j \qquad (4)$$

$$\forall l, \ \xi_l \geq 0.0 \qquad (5)$$

$$\forall \text{ random switch } \Xi, \ \sum_{l : \xi_l \in \Xi} \xi_l = 1.0 \qquad (6)$$

Equation 2 computes the limiting distribution, $\mathbf{y}$, of the embedded Discrete Time Markov Chain. Furthermore, the steady-state probabilities, $\pi = [\pi_k]$, for the underlying continuous-time stochastic process, are obtained by Equation 3. In Equation 3, $E[s_k]$ denotes the expected sojourn time for tangible marking $s_k \in \mathcal{T}$, and it is computed by Equation 4, where $r_j$ denotes the (firing) rate of (timed) transition $T_j$. Equations 5 and 6, introduce the notion of dynamic random switch in the formulation, constraining appropriately the decision variables $\xi_l$. Finally, once the steady-state probability vector $\pi$ has been obtained, various performance measures of interest can be defined as appropriate functions of $\pi$ and the other system parameters. For instance, Equation 1 can represent the RAS (steady-state) throughput function, by setting

$$\mathcal{O}(\xi) \equiv \sum_{(k,j)} \pi_k r_j I_{\{T_j \text{ enabled in } s_k \\ \land \text{ cor. to a job} \\ \text{ unloading}\}} \qquad (7)$$

## 3. DETERMINISTIC OPTIMAL SCHEDULING POLICIES

It can be shown that the solution space of the MP formulation employing the objective function defined by Equation 7, will always contain a *deterministic* optimal solution, i.e., an optimal solution $\xi^*$ with $\xi_l^* \in \{0, 1\}$, $\forall l$. This result can be obtained through direct analysis of the structure of the formulation under consideration (Choi and Reveliotis, 2001), and it is consistent to some more classical results in Dynamic Programming theory (Puterman, 1994). Beyond its theoretical interest, establishing the existence of a deterministic optimal solution is of practical importance because it limits the search for an optimal solution over a discrete sub-space of the overall solution space, and therefore, it renders the considered problem amenable to enumerative approaches. We notice, however, that the aforementioned deterministic structure of the optimal solution can be lost, once some additional constraint(s) is introduced in the formulation. For instance, in a multi-item production scheduling context, such a constraint can arise from the requirement for production of the various product types at quantities that satisfy certain ratios; we refer the reader to (Choi and Reveliotis, 2001) for a more concrete example.

## 4. AN EXAMPLE: THROUGHPUT MAXIMIZATION OF A CAPACITATED RE-ENTRANT LINE

**The considered RAS** This section, elucidates the theory developed in Section 2 and 3, by applying the developed results to the detailed analysis of the small capacitated re-entrant line in Figure 2. This line produces a single item, and it possesses two stations, $W_1$ and $W_2$, with $S_1 = S_2 = 1$ and $C_1 = 1$; $C_2 = 2$, where $S_i$ is the number of identical servers and $C_i$ the number of buffer slots in the workstation $W_i$. Furthermore, the supported production sequence is $J = < J_1, J_2, J_3 >$, with $W(J_1) = W(J_3) = W_1$ and $W(J_2) = W_2$. Finally, stage processing times are exponentially distributed with means $m_j = 1/\mu_j > 0$, $j = 1, 2, 3$, and so are the involved transfer times, with a uniform mean $d = 1/\lambda$. Each part visiting a workstation for the execution of some processing stage is allocated one unit of buffering capacity, which it holds exclusively during its entire sojourn in the station. Once in the station local buffer, the part competes for one of the station servers, for the execution of the requested stage. A part having finished the processing of its current stage at a certain station, waits in its allocated buffer for transfer to the next requested station. This transfer is facilitated by the central (automated) material handling system, and it is authorized
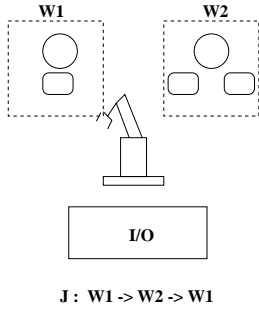
W1   W2

**J : W1 -> W2 -> W1**

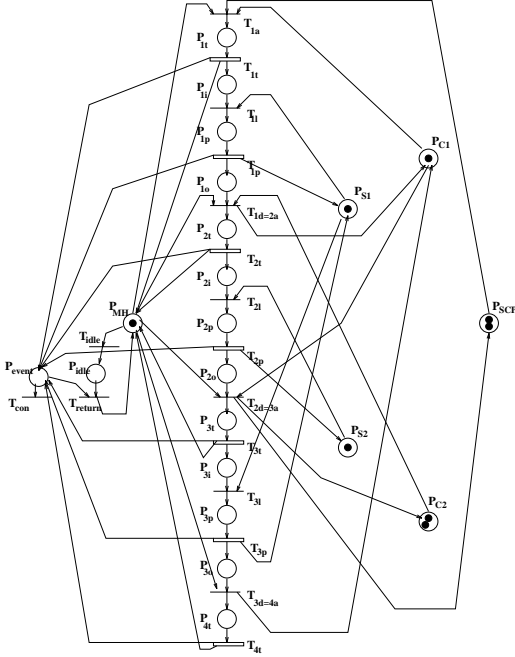Fig. 1. Example: The capacitated re-entrant line



Fig. 2. Example: The GSPN model

by a supervisory control policy ensuring that (i) the destination workstation has available buffering capacity, and (ii) the transfer is *safe*, i.e., it is still physically possible from the resulting state to process all running jobs to completion. For this small configuration, it is easy to see that, under the operational assumptions outlined above, the system material flow will remain deadlock-free, as long as

$$|J_1| + |J_2| \leq C_1 + C_2 - 1 = 2 \qquad (8)$$

where $|J_j|$, $j = 1, 2, 3$ denotes the number of job instances in $W(J_j)$ executing stage $J_j$.

**GSPN-based modeling of the capacitated re-entrant line** The GSPN modeling the behavior of the capacitated re-entrant line of Figure 2, under the control of the maximally permissive structural control policy (SCP) of Equation 8, is depicted in Figure 3. Specifically, in the GSPN of Figure 3, the part flow dynamics associated with each processing stage $J_j$, $j = 1, 2, 3$, are modeled by the corresponding net path $< T_{ja}, P_{jt}, T_{jt}, P_{ji}, T_{jl}, P_{jp}, T_{jp}, P_{jo}, T_{jd} >$, while it also holds $T_{jd} \equiv T_{j+1,a}$, with $j = 4$ denot-

ing the last unloading step. A token in place $P_{jt}$ represents a part in transit to the buffer of workstation $W(J_j)$; a token in place $P_{ji}$ represents a part in the buffer of $W(J_j)$ waiting the allocation of one of the buffer servers; a token in place $P_{jp}$ represents a part in processing of stage $J_j$; finally, a token in place $P_{jo}$ represents a part having finished processing of stage $J_j$, and waiting for transfer to the next requested workstation or, in case that $J_j$ is the last processing stage, to the I/O station. On the other hand, places $P_{MH}$, $P_{S_i}$, $P_{C_i}$, $i = 1, 2$, and $P_{SCP}$ model respectively the availability of the system transporter, workstation servers and buffers, and the logic of the applied SCP, according to the standard, by now, modeling practice of resource-process nets (Banaszak and Krogh, 1990). It is important to notice that transitions $T_{ja}$, $T_{jl}$ and $T_{jd}$, that are associated with the various decisions regarding the allocation of the system buffering, processing and/or transport capacity, are untimed/ immediate transitions, while the delays experienced from the processing and/or transfer times involved with the execution of these decisions, are modeled by the timed transitions $T_{jt}$ and $T_{jp}$. This separation of the net components modeling the timings of the various system events from the net structure modeling the underlying resource allocation and the associated decision making, enables the modeling of timing distributions other than exponential through the (local) substitution of the corresponding timed transitions by GSPN subnets modeling the approximating phase-type distributions (c.f. Section 2). It also allows the modeling of the required scheduling logic through a set of *dynamic random switches*, that resolve the conflicts among the immediate transitions that are simultaneously enabled at the net reachable vanishing markings. Finally, some explanation is necessary about the role of places $P_{idle}$, $P_{event}$ and their associated transitions $T_{idle}$, $T_{return}$ and $T_{con}$. This subnet essentially establishes a GSPN-compatible mechanism for representing some deliberate idleness in the underlying scheduling logic, since, in the considered operational context, the optimal scheduling policy is not necessarily non-idling. Hence, the triggering of transition $T_{idle}$ consumes the transporter-modeling token, which remains in place $P_{idle}$, until the immediate transition $T_{return}$ is enabled through the presence of a token in place $P_{event}$. $P_{event}$ is marked every time that one of the system timed transitions fires, signaling the completion of some event. Notice that $T_{return}$ will always be in conflict with transition $T_{con}$, but it is assumed to have priority over the latter, which is technically imposed by setting the corresponding (static) random switch to $\{\xi_{T_{return}} = 1, \xi_{T_{con}} = 0\}$. Finally, $T_{con}$ is a sink transition that "consumes" event completion signaling tokens, in case that the transporter is not (deliberately) idling.

Fig. 3. Example: The Embedded Markov Chain (EMC)

**Maximizing the line throughput** The EMC for the GSPN of Figure 3 is presented in Figure 4, while the net markings corresponding to the various states depicted in Figure 4 are listed in Table 1, at the end of the document. In Figure 4, states corresponding to vanishing markings are depicted by single circles, while states corresponding to tangible markings are depicted by double circles. Furthermore, the part of the chain depicted in dashed lines should be inaccessible under operation by any optimal scheduling policy, either because it leads to dead/absorbing states [1] – e.g., transition from $s_6$ to $s_7$ – or because the transitions branching to that part of the chain essentially introduce some unnecessary delay in the system operation, by deliberately idling the server – e.g., transition from $s_{30}$ to $s_{31}$. The remaining modified EMC, depicted with solid lines in Figure 4, contains only two random switches of two options each, which combined with Equation 6, leaves us with two decision variables $\xi_1$ and $\xi_2$. Then, the structure of the optimal scheduling policy can be obtained by computing the closed-form expressions for $TH(0,0)$, $TH(0,1)$, $TH(1,0)$ and $TH(1,1)$, by means of Equations 2 – 4 and 7, and determining the parameter ranges over which each of these expressions dominates the others. Working according to this plan, one can establish that the dominance relationships among these four expressions are those depicted by the lattice of Figure 5. The reader can verify that the optimal policy, defined by $(\xi_1 = 1, \xi_2 = 1)$, essentially implements a First-Buffer-First-Serve (FBFS) logic in the considered operational context.

---

[1] These absorbing states are due to the mechanism that introduces deliberate system idleness, discussed above: specifically, the firing of transition $T_{idle}$ in a vanishing marking with no pending timed events can trap for ever the token modeling the material-handling availability in place $P_{idle}$.



Fig. 4. Example: Characterizing the dominance among the candidate scheduling policies

## 5. DISCUSSION

While the proposed approach provides a detailed analytical characterization of the performance control problem for structurally controlled RAS, from an implementational standpoint, it requires the explicit enumeration of the underlying state space, which explodes very fast. This limits the applicability of the presented methodology as a practical scheduling approach, and necessitates the development of pertinent approximating schemes, that will lead to (near-)optimal scheduling policies for structurally controlled RAS. The development of such effective approximating schemes is part of our current investigations.

## ACKNOWLEDGMENT

## REFERENCES

A. Marsan, M., G. Balbo and G. Conte (1986). *Performance Models of Multiprocessor Systems.* The MIT Press. Cambridge, MA.

Banaszak, Z. A. and B. H. Krogh (1990). Deadlock avoidance in flexible manufacturing systems with concurrently competing process flows. *IEEE Trans. on Robotics and Automation* **6**, 724–734.

Choi, J. Y. and S. Reveliotis (2001). Performance modeling and control of capacitated re-entrant lines. Technical report. ISyE, Georgia Institute of Technology.

Gershwin, S. B. (1989). Hierarchical flow control: A framework for scheduling and planning discrete events in manufacturing systems. *Proceedings of the IEEE* **77**, 195–209.

Papadopoulos, H. T., C. Heavy and J. Browne (1993). *Queueing Theory in Manufacturing Systems Analysis and Design.* Chapman & Hall. New York, NY.

Puterman, M. L. (1994). *Markov Decision Processes: Discrete Stochastic Dynamic Programming.* John Wiley & Sons.

Reveliotis, S. (2001). Liveness enforcing supervision for sequential resource allocation systems: State of the art and open issues. In:

Table 1. Example:The EMC markings

| $s_k$ | $P_{1t}P_{1i}P_{1p}P_{1o}$ | $P_{2t}P_{2i}P_{2p}P_{2o}$ | $P_{3t}P_{3i}P_{3p}P_{3o}P_{4t}$ | $P_{MH}P_{idle}P_{event}$ | $P_{S_1}P_{S_2}$ | $P_{C_1}P_{C_2}P_{SCP}$ |
|---|---|---|---|---|---|---|
| 0 | 0 0 0 0 | 0 0 0 0 | 0 0 0 0 0 | 1 0 0 | 1 1 | 1 2 2 |
| 1 | 1 0 0 0 | 0 0 0 0 | 0 0 0 0 0 | 0 0 0 | 1 1 | 0 2 1 |
| 2 | 0 0 0 0 | 0 0 0 0 | 0 0 0 0 0 | 0 1 0 | 1 1 | 1 2 2 |
| 3 | 0 1 0 0 | 0 0 0 0 | 0 0 0 0 0 | 1 0 1 | 1 1 | 0 2 1 |
| 4 | 0 0 1 0 | 0 0 0 0 | 0 0 0 0 0 | 0 1 0 | 0 1 | 0 2 1 |
| 5 | 0 0 0 1 | 0 0 0 0 | 0 0 0 0 0 | 0 1 1 | 1 1 | 0 2 1 |
| 6 | 0 0 0 1 | 0 0 0 0 | 0 0 0 0 0 | 1 0 0 | 1 1 | 0 2 1 |
| 7 | 0 0 0 1 | 0 0 0 0 | 0 0 0 0 0 | 0 1 0 | 1 1 | 0 2 1 |
| 8 | 0 0 0 0 | 1 0 0 0 | 0 0 0 0 0 | 0 0 0 | 1 1 | 1 1 1 |
| 9 | 0 0 0 0 | 0 1 0 0 | 0 0 0 0 0 | 1 0 1 | 1 1 | 1 1 1 |
| 10 | 0 0 0 0 | 0 0 1 0 | 0 0 0 0 0 | 1 0 0 | 1 0 | 1 1 1 |
| 11 | 0 0 0 0 | 0 0 1 0 | 0 0 0 0 0 | 0 1 0 | 1 0 | 1 1 1 |
| 12 | 1 0 0 0 | 0 0 1 0 | 0 0 0 0 0 | 0 0 0 | 1 0 | 0 1 0 |
| 13 | 0 0 0 0 | 0 0 0 1 | 0 0 0 0 0 | 0 1 1 | 1 1 | 1 1 1 |
| 14 | 0 0 0 0 | 0 0 0 1 | 0 0 0 0 0 | 1 0 0 | 1 1 | 1 1 1 |
| 15 | 0 0 0 0 | 0 0 0 1 | 0 0 0 0 0 | 0 1 0 | 1 1 | 1 1 1 |
| 16 | 0 0 0 0 | 0 0 0 0 | 1 0 0 0 0 | 0 0 0 | 1 1 | 0 2 2 |
| 17 | 1 0 0 0 | 0 0 0 1 | 0 0 0 0 0 | 0 0 0 | 1 1 | 0 1 0 |
| 18 | 0 0 0 0 | 0 0 0 0 | 0 1 0 0 0 | 1 0 1 | 1 1 | 0 2 2 |
| 19 | 0 0 0 0 | 0 0 0 0 | 0 0 1 0 0 | 0 1 0 | 0 1 | 0 2 2 |
| 20 | 0 0 0 0 | 0 0 0 0 | 0 0 0 1 0 | 0 1 1 | 1 1 | 0 2 2 |
| 21 | 0 0 0 0 | 0 0 0 0 | 0 0 0 1 0 | 1 0 0 | 1 1 | 0 2 2 |
| 22 | 0 0 0 0 | 0 0 0 0 | 0 0 0 1 0 | 0 1 0 | 1 1 | 0 2 2 |
| 23 | 0 0 0 0 | 0 0 0 0 | 0 0 0 0 1 | 0 0 0 | 1 1 | 1 2 2 |
| 24 | 0 0 0 0 | 0 0 0 0 | 0 0 0 0 0 | 1 0 1 | 1 1 | 1 2 2 |
| 25 | 0 1 0 0 | 0 0 1 0 | 0 0 0 0 0 | 1 0 1 | 1 0 | 0 1 0 |
| 26 | 1 0 0 0 | 0 0 0 1 | 0 0 0 0 0 | 0 0 1 | 1 1 | 0 1 0 |
| 27 | 0 0 1 0 | 0 0 1 0 | 0 0 0 0 0 | 0 1 0 | 0 0 | 0 1 0 |
| 28 | 0 0 0 1 | 0 0 1 0 | 0 0 0 0 0 | 0 1 1 | 1 0 | 0 1 0 |
| 29 | 0 0 1 0 | 0 0 0 1 | 0 0 0 0 0 | 0 1 1 | 0 1 | 0 1 0 |
| 30 | 0 0 0 1 | 0 0 1 0 | 0 0 0 0 0 | 1 0 0 | 1 0 | 0 1 0 |
| 31 | 0 0 0 1 | 0 0 1 0 | 0 0 0 0 0 | 0 1 0 | 1 0 | 0 1 0 |
| 32 | 0 0 0 0 | 1 0 1 0 | 0 0 0 0 0 | 0 0 0 | 1 0 | 1 0 0 |
| 33 | 0 0 0 1 | 0 0 0 1 | 0 0 0 0 0 | 0 1 1 | 1 1 | 0 1 0 |
| 34 | 0 0 0 1 | 0 0 0 1 | 0 0 0 0 0 | 1 0 0 | 1 1 | 0 1 0 |
| 35 | 0 0 0 0 | 1 0 0 1 | 0 0 0 0 0 | 0 0 0 | 1 1 | 1 0 0 |
| 36 | 0 0 0 1 | 0 0 0 1 | 0 0 0 0 0 | 0 1 0 | 1 1 | 0 1 0 |
| 37 | 0 0 0 0 | 0 1 0 1 | 0 0 0 0 0 | 1 0 1 | 1 1 | 1 0 0 |
| 38 | 0 0 0 0 | 0 0 1 1 | 0 0 0 0 0 | 0 1 0 | 1 0 | 1 0 0 |
| 39 | 0 0 0 0 | 0 0 0 2 | 0 0 0 0 0 | 0 1 1 | 1 1 | 1 0 0 |
| 40 | 0 0 0 0 | 0 0 0 2 | 0 0 0 0 0 | 1 0 0 | 1 1 | 1 0 0 |
| 41 | 0 0 0 0 | 0 0 0 2 | 0 0 0 0 0 | 0 1 0 | 1 1 | 1 0 0 |
| 42 | 0 0 0 0 | 0 0 0 1 | 1 0 0 0 0 | 0 0 0 | 1 1 | 0 1 1 |
| 43 | 0 0 0 0 | 0 0 0 1 | 0 1 0 0 0 | 1 0 1 | 1 1 | 0 1 1 |
| 44 | 0 0 0 0 | 0 0 0 1 | 0 0 1 0 0 | 0 1 0 | 0 1 | 0 1 1 |
| 45 | 0 0 0 0 | 0 0 0 1 | 0 0 0 1 0 | 0 1 1 | 1 1 | 0 1 1 |
| 46 | 0 0 0 0 | 0 0 0 1 | 0 0 0 1 0 | 1 0 0 | 1 1 | 0 1 1 |
| 47 | 0 0 0 0 | 0 0 0 1 | 0 0 0 1 0 | 0 1 0 | 1 1 | 0 1 1 |
| 48 | 0 0 0 0 | 0 0 0 1 | 0 0 0 0 1 | 0 0 0 | 1 1 | 1 1 1 |
| 49 | 0 0 0 0 | 0 0 0 1 | 0 0 0 0 0 | 1 0 1 | 1 1 | 1 1 1 |
| 50 | 0 1 0 0 | 0 0 0 1 | 0 0 0 0 0 | 1 0 1 | 1 1 | 0 1 0 |
| 51 | 0 0 1 0 | 0 0 0 1 | 0 0 0 0 0 | 0 1 0 | 0 1 | 0 1 0 |
| 52 | 0 0 1 0 | 0 0 0 1 | 0 0 0 0 0 | 1 0 0 | 0 1 | 0 1 0 |
| 53 | 0 0 0 0 | 0 1 1 0 | 0 0 0 0 0 | 1 0 1 | 1 0 | 1 0 0 |
| 54 | 0 0 0 0 | 1 0 0 1 | 0 0 0 0 0 | 0 0 1 | 1 1 | 1 0 0 |
| 55 | 0 0 0 0 | 0 1 1 0 | 0 0 0 0 0 | 0 1 0 | 1 0 | 1 0 0 |
| 56 | 0 0 0 0 | 0 1 0 1 | 0 0 0 0 0 | 0 1 1 | 1 1 | 1 0 0 |
| 57 | 0 0 0 0 | 0 0 1 1 | 0 0 0 0 0 | 1 0 0 | 1 0 | 1 0 0 |
| 58 | 0 0 0 0 | 0 0 1 0 | 1 0 0 0 0 | 0 0 0 | 1 0 | 0 1 1 |
| 59 | 0 0 0 0 | 0 0 0 1 | 1 0 0 0 0 | 0 0 1 | 1 1 | 0 1 1 |
| 60 | 0 0 0 0 | 0 0 1 0 | 0 1 0 0 0 | 1 0 1 | 1 0 | 0 1 1 |
| 61 | 0 0 0 0 | 0 0 1 0 | 0 0 1 0 0 | 0 1 0 | 0 0 | 0 1 1 |
| 62 | 0 0 0 0 | 0 0 0 1 | 0 0 1 0 0 | 0 1 1 | 0 1 | 0 1 1 |
| 63 | 0 0 0 0 | 0 0 1 0 | 0 0 0 1 0 | 0 1 1 | 1 0 | 0 1 1 |
| 64 | 0 0 0 0 | 0 0 1 0 | 0 0 0 1 0 | 1 0 0 | 1 0 | 0 1 1 |
| 65 | 0 0 0 0 | 0 0 1 0 | 0 0 0 1 0 | 0 1 0 | 1 0 | 0 1 1 |
| 66 | 0 0 0 0 | 0 0 1 0 | 0 0 0 0 1 | 0 0 0 | 1 0 | 1 1 1 |
| 67 | 0 0 0 0 | 0 0 0 1 | 0 0 0 1 0 | 0 1 1 | 1 1 | 0 1 1 |
| 68 | 0 0 0 0 | 0 0 1 0 | 0 0 0 0 0 | 1 0 1 | 1 0 | 1 1 1 |
| 69 | 0 0 0 0 | 0 0 0 1 | 0 0 1 0 0 | 1 0 0 | 0 1 | 0 1 1 |
| 70 | 0 0 0 0 | 0 0 0 1 | 0 0 0 0 1 | 0 0 1 | 1 1 | 1 1 1 |