

Datapath Subsystems

11

11.1 Introduction

Chip functions generally can be divided into the following categories:

- Datapath operators
- Memory elements
- Control structures
- Special-purpose cells
 - I/O
 - Power distribution
 - Clock generation and distribution
 - Analog and RF

CMOS system design consists of partitioning the system into subsystems of the types listed above. Many options exist that make trade-offs between speed, density, programmability, ease of design, and other variables. This chapter addresses design options for common datapath operators. The next chapter addresses arrays, especially those used for memory. Control structures are most commonly coded in a hardware description language and synthesized. Special-purpose subsystems are considered in Chapter 13.

As introduced in Chapter 1, datapath operators benefit from the structured design principles of hierarchy, regularity, modularity, and locality. They may use N identical circuits to process N -bit data. Related data operators are placed physically adjacent to each other to reduce wire length and delay. Generally, data is arranged to flow in one direction, while control signals are introduced in a direction orthogonal to the dataflow.

Common datapath operators considered in this chapter include adders, one/zero detectors, comparators, counters, Boolean logic units, error-correcting code blocks, shifters, and multipliers.

11.2 Addition/Subtraction

“Multitudes of contrivances were designed, and almost endless drawings made, for the purpose of economizing the time and simplifying the mechanism of carriage.”

—Charles Babbage, on Difference Engine No. 1, 1864 [Morrison61]

Addition forms the basis for many processing operations, from ALUs to address generation to multiplication to filtering. As a result, adder circuits that add two binary numbers are of great interest to digital system designers. An extensive, almost endless, assortment of adder architectures serve different speed/power/area requirements. This section begins with half adders and full adders for single-bit addition. It then considers a plethora of carry-propagate adders (CPAs) for the addition of multibit words. Finally, related structures such as subtracters and multiple-input adders are discussed.

11.2.1 Single-Bit Addition

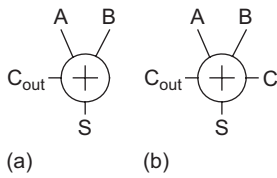


FIGURE 11.1
Half and full adders

The *half adder* of Figure 11.1(a) adds two single-bit inputs, A and B . The result is 0, 1, or 2, so two bits are required to represent the value; they are called the sum S and carry-out C_{out} . The carry-out is equivalent to a carry-in to the next more significant column of a multibit adder, so it can be described as having double the *weight* of the other bits. If multiple adders are to be cascaded, each must be able to receive the carry-in. Such a *full adder* as shown in Figure 11.1(b) has a third input called C or C_{in} .

The truth tables for the half adder and full adder are given in Tables 11.1 and 11.2. For a full adder, it is sometimes useful to define *Generate* (G), *Propagate* (P), and *Kill* (K) signals. The adder generates a carry when C_{out} is true independent of C_{in} , so $G = A \cdot B$. The adder kills a carry when C_{out} is false independent of C_{in} , so $K = \overline{A} \cdot \overline{B} = \overline{A + B}$. The adder propagates a carry; i.e., it produces a carry-out if and only if it receives a carry-in, when exactly one input is true: $P = A \oplus B$.

TABLE 11.1 Truth table for half adder

A	B	C_{out}	S
0	0	0	0
0	1	0	1
1	0	0	1
1	1	1	0

TABLE 11.2 Truth table for full adder

A	B	C	G	P	K	C_{out}	S
0	0	0	0	0	1	0	0
		1				0	1
0	1	0	0	1	0	0	1
		1				1	0
1	0	0	0	1	0	0	1
		1				1	0
1	1	0	1	0	0	1	0
		1				1	1

From the truth table, the half adder logic is

$$\begin{aligned} S &= A \oplus B \\ C_{out} &= A \cdot B \end{aligned} \quad (11.1)$$

and the full adder logic is

$$\begin{aligned}
 S &= \overline{A}\overline{B}C + \overline{A}B\overline{C} + A\overline{B}\overline{C} + ABC \\
 &= (A \oplus B) \oplus C = P \oplus C \\
 C_{\text{out}} &= AB + AC + BC \\
 &= AB + C(A + B) \\
 &= \overline{\overline{AB} + \overline{C}(\overline{A} + \overline{B})} \\
 &= \text{MAJ}(A, B, C)
 \end{aligned} \tag{11.2}$$

The most straightforward approach to designing an adder is with logic gates. Figure 11.2 shows a half adder. Figure 11.3 shows a full adder at the gate (a) and transistor (b) levels. The carry gate is also called a *majority* gate because it produces a 1 if at least two of the three inputs are 1. Full adders are used most often, so they will receive the attention of the remainder of this section.

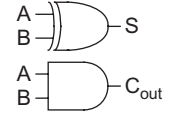


FIGURE 11.2
Half adder design

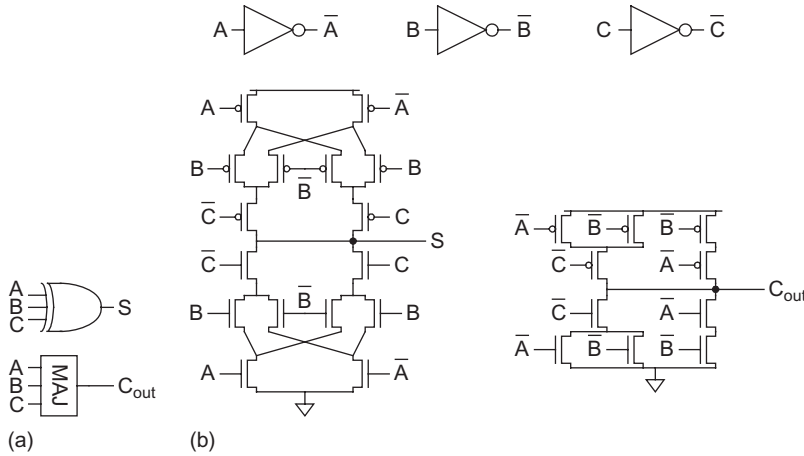


FIGURE 11.3 Full adder design

The full adder of Figure 11.3(b) employs 32 transistors (6 for the inverters, 10 for the majority gate, and 16 for the 3-input XOR). A more compact design is based on the observation that S can be factored to reuse the C_{out} term as follows:

$$S = ABC + (A + B + C)\overline{C}_{\text{out}} \tag{11.3}$$

Such a design is shown at the gate (a) and transistor (b) levels in Figure 11.4 and uses only 28 transistors. Note that the pMOS network is identical to the nMOS network rather than being the conduction complement, so the topology is called a *mirror adder*. This simplification reduces the number of series transistors and makes the layout more uniform. It is possible because the addition function is *symmetric*; i.e., the function of complemented inputs is the complement of the function.

The mirror adder has a greater delay to compute S than C_{out} . In carry-ripple adders (Section 11.2.2.1), the critical path goes from C to C_{out} through many full adders, so the

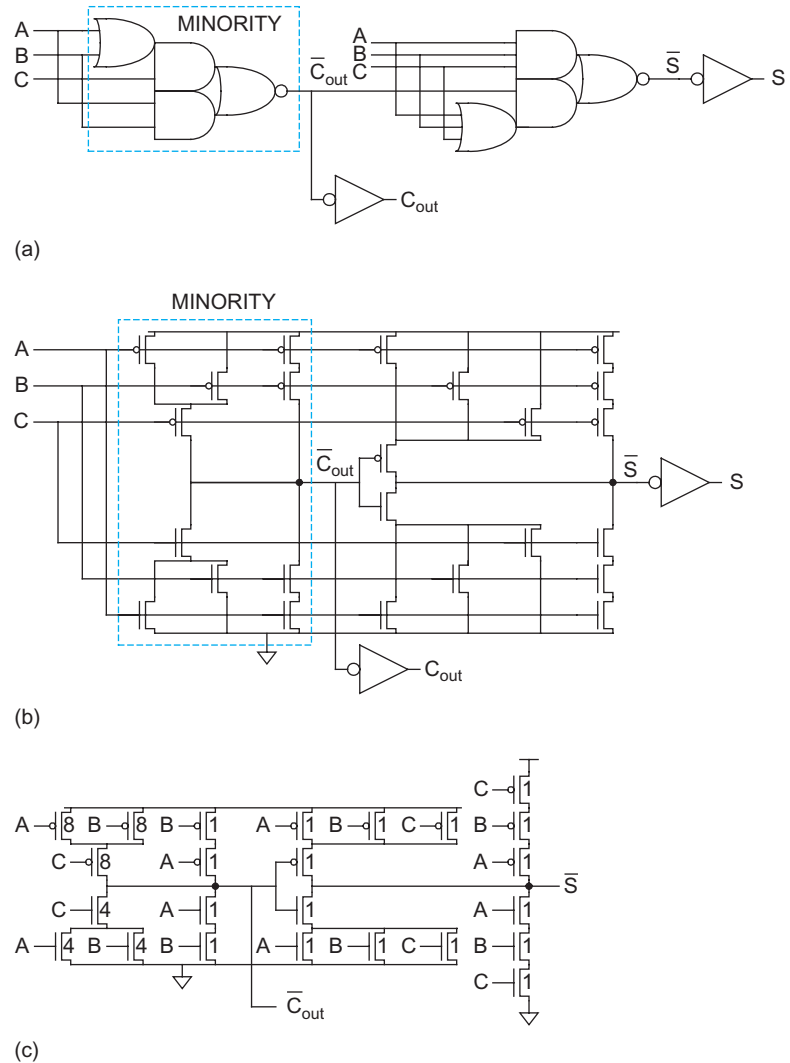


FIGURE 11.4 Full adder for carry-ripple operation

extra delay computing S is unimportant. Figure 11.4(c) shows the adder with transistor sizes optimized to favor the critical path using a number of techniques:

- Feed the carry-in signal (C) to the inner inputs so the internal capacitance is already discharged.
- Make all transistors in the sum logic whose gate signals are connected to the carry-in and carry logic minimum size (1 unit, e.g., 4λ). This minimizes the branching effort on the critical path. Keep routing on this signal as short as possible to reduce interconnect capacitance.
- Determine widths of series transistors by logical effort and simulation. Build an asymmetric gate that reduces the logical effort from C to C_{out} at the expense of effort to S .

- Use relatively large transistors on the critical path so that stray wiring capacitance is a small fraction of the overall capacitance.
- Remove the output inverters and alternate positive and negative logic to reduce delay and transistor count to 24 (see Section 11.2.2.1).

Figure 11.5 shows two layouts of the adder (see also the inside front cover). The choice of the aspect ratio depends on the application. In a standard-cell environment, the layout of Figure 11.5(a) might be appropriate when a single row of nMOS and pMOS transistors is used. The routing for the A , B , and C inputs is shown inside the cell, although it could be placed outside the cell because external routing tracks have to be assigned to these signals anyway. Figure 11.5(b) shows a layout that might be appropriate for a dense datapath (if horizontal polysilicon is legal). Here, the transistors are rotated and all of the wiring is completed in polysilicon and metal1. This allows metal2 bus lines to pass over the cell horizontally. Moreover, the widths of the transistors can increase

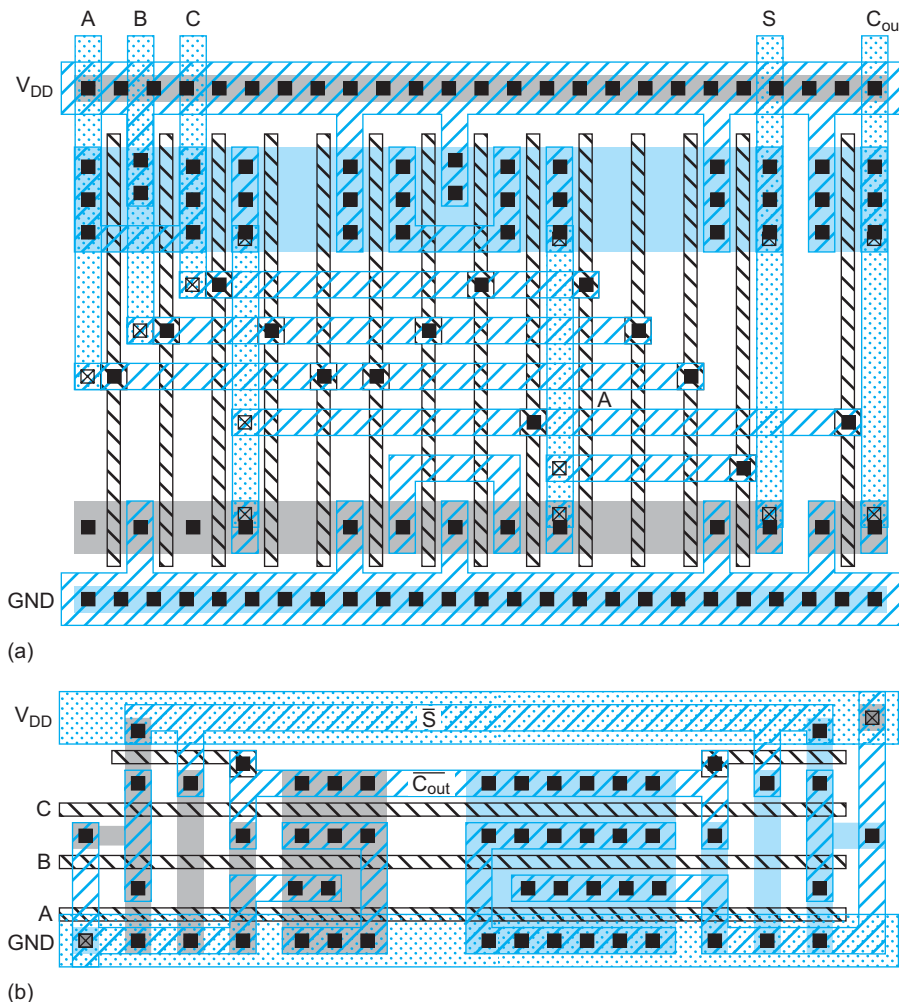


FIGURE 11.5 Full adder layouts. Color version on inside front cover.

without impacting the bit-pitch (height) of the datapath. In this case, the widths are selected to reduce the C_{in} to C_{out} delay that is on the critical path of a carry-ripple adder.

A rather different full adder design uses transmission gates to form multiplexers and XORs. Figure 11.6(a) shows the transistor-level schematic using 24 transistors and providing buffered outputs of the proper polarity with equal delay. The design can be understood by parsing the transmission gate structures into multiplexers and an “invertible inverter” XOR structure (see Section 11.7.4), as drawn in Figure 11.6(b).¹ Note that the multiplexer choosing S is configured to compute $P \oplus C$, as given in EQ (11.2).

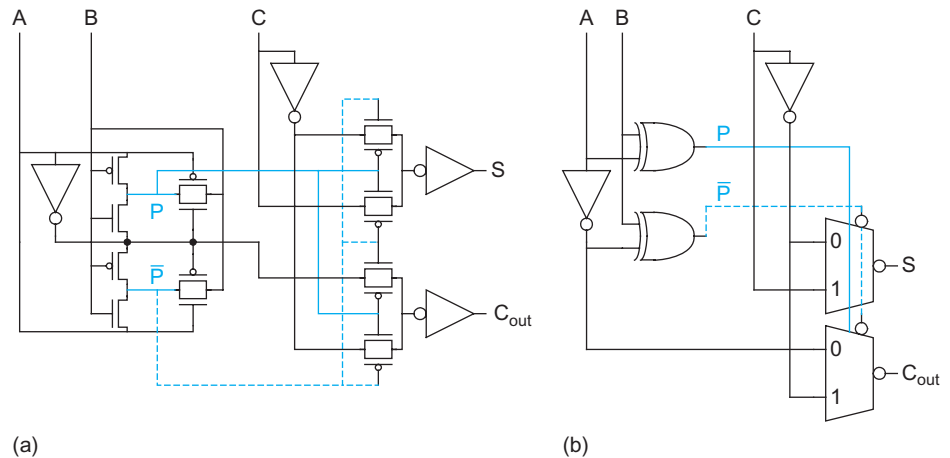


FIGURE 11.6 Transmission gate full adder

Figure 11.7 shows a complementary pass-transistor logic (CPL) approach. In comparison to a poorly optimized 40-transistor static CMOS full adder, [Yano90] finds CPL is twice as fast, 30% lower in power, and slightly smaller. On the other hand, in comparison to a careful implementation of the mirror adder, [Zimmermann97] finds the CPL delay slightly better, the power comparable, and the area much larger.

Dynamic full adders are widely used in fast multipliers when power is not a concern. As the sum logic inherently requires true and complementary versions of the inputs, dual-rail domino is necessary. Figure 11.8 shows such an adder using footless dual-rail domino XOR/XNOR and MAJORITY/MINORITY gates [Heikes94]. The delays to the two outputs are reasonably well balanced, which is important for multipliers where both paths are critical. It shares transistors in the sum gate to reduce transistor count and takes advantage of the symmetric property to provide identical layouts for the two carry gates.

Static CMOS full adders typically have a delay of 2–3 FO4 inverters, while domino adders have a delay of about 1.5.

11.2.2 Carry-Propagate Addition

N -bit adders take inputs $\{A_N, \dots, A_1\}$, $\{B_N, \dots, B_1\}$, and carry-in C_{in} , and compute the sum $\{S_N, \dots, S_1\}$ and the carry-out of the most significant bit C_{out} , as shown in Figure 11.9.

¹Some switch-level simulators, notably IRSIM, are confused by this XOR structure and may not simulate it correctly.

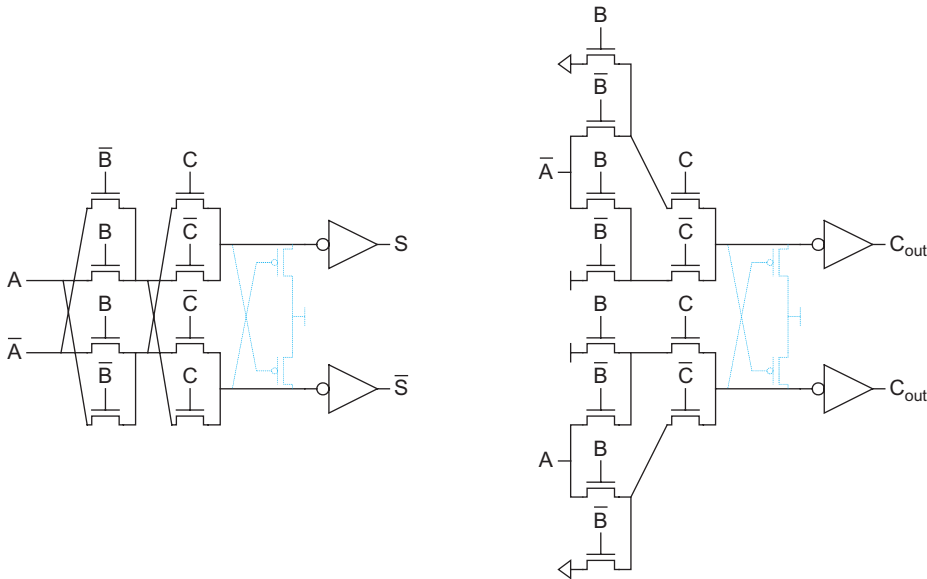


FIGURE 11.7 CPL full adder

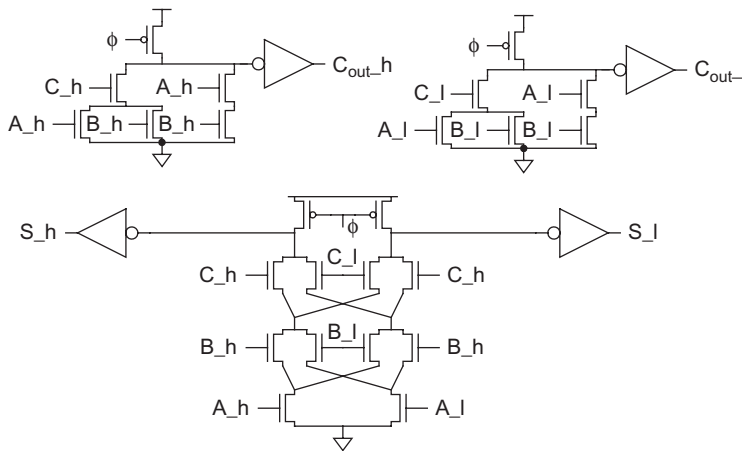


FIGURE 11.8 Dual-rail domino full

(Ordinarily, this text calls the least significant bit A_0 rather than A_1 . However, for adders, the notation developed on subsequent pages is more graceful if column 0 is reserved to handle the carry.) They are called *carry-propagate adders* (CPAs) because the carry into each bit can influence the carry into all subsequent bits. For example, Figure 11.10 shows the addition $1111_2 + 0000_2 + 0/1$, in which each of the sum and carry bits is influenced by C_{in} . The simplest design is the carry-ripple adder in which the carry-out of one bit is simply connected as the carry-in to the next. Faster adders look ahead to predict the carry-out of a multibit group. This is usually done by computing group PG signals to indicate

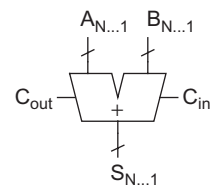


FIGURE 11.9 Carry-propagate adder

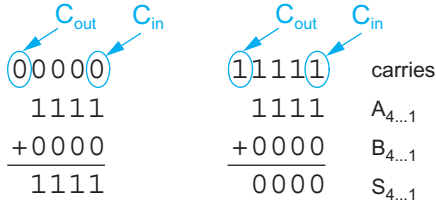


FIGURE 11.10 Example of carry propagation

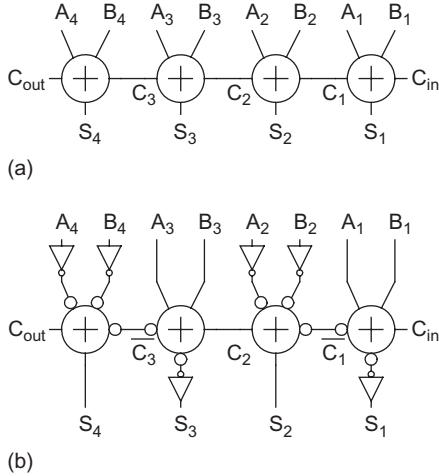


FIGURE 11.11 4-bit carry-ripple adder

whether the multibit group will propagate a carry-in or will generate a carry-out. Long adders use multiple levels of lookahead structures for even more speed.

11.2.2.1 Carry-Ripple Adder An N -bit adder can be constructed by cascading N full adders, as shown in Figure 11.11(a) for $N = 4$. This is called a *carry-ripple adder* (or *ripple-carry adder*). The carry-out of bit i , C_i , is the carry-in to bit $i + 1$. This carry is said to have twice the *weight* of the sum S_i . The delay of the adder is set by the time for the carries to ripple through the N stages, so the $t_{C \rightarrow C_{out}}$ delay should be minimized.

This delay can be reduced by omitting the inverters on the outputs, as was done in Figure 11.4(c). Because addition is a *self-dual function* (i.e., the function of complementary inputs is the complement of the function), an inverting full adder receiving complementary inputs produces true outputs. Figure 11.11(b) shows a carry-ripple adder built from inverting full adders. Every other stage operates on complementary data. The delay inverting the adder inputs or sum outputs is off the critical ripple-carry path.

11.2.2.2 Carry Generation and Propagation This section introduces notation commonly used in describing faster adders. Recall that the P (*propagate*) and G (*generate*) signals were defined in Section 11.2.1. We can generalize these signals to describe whether a group spanning bits $i \dots j$, inclusive, generate a carry or propagate a carry. A group of bits generates a carry if its carry-out is true independent of the carry-in; it propagates a carry if its carry-out is true when there is a carry-in. These signals can be defined recursively for $i \geq k > j$ as

$$\begin{aligned} G_{i:j} &= G_{i:k} + P_{i:k} \cdot G_{k-1:j} \\ P_{i:j} &= P_{i:k} \cdot P_{k-1:j} \end{aligned} \quad (11.4)$$

with the base case

$$\begin{aligned} G_{i:i} &\equiv G_i = A_i \cdot B_i \\ P_{i:i} &\equiv P_i = A_i \oplus B_i \end{aligned} \quad (11.5)$$

In other words, a group generates a carry if the upper (more significant) or the lower portion generates and the upper portion propagates that carry. The group propagates a carry if both the upper and lower portions propagate the carry.²

The carry-in must be treated specially. Let us define $C_0 = C_{in}$ and $C_N = C_{out}$. Then we can define generate and propagate signals for bit 0 as

$$\begin{aligned} G_{0:0} &= C_{in} \\ P_{0:0} &= 0 \end{aligned} \quad (11.6)$$

²Alternatively, many adders use $\bar{K}_i = A_i + B_i$ in place of P_i because OR is faster than XOR. The group logic uses the same gates: $G_{i:j} = G_{i:k} + \bar{K}_{i:k} \cdot G_{k-1:j}$ and $\bar{K}_{i:j} = \bar{K}_{i:k} \cdot \bar{K}_{k-1:j}$. However, $P_i = A_i \oplus B_i$ is still required in EQ (11.7) to compute the final sum. It is sometimes renamed X_i or T_i to avoid ambiguity.

Observe that the carry into bit i is the carry-out of bit $i-1$ and is $C_{i-1} = G_{i-1:0}$. This is an important relationship; *group generate* signals and *carries* will be used synonymously in the subsequent sections. We can thus compute the sum for bit i using EQ(11.2) as

$$S_i = P_i \oplus G_{i-1:0} \quad (11.7)$$

Hence, addition can be reduced to a three-step process:

1. Computing bitwise generate and propagate signals using EQs (11.5) and (11.6)
2. Combining PG signals to determine group generates $G_{i-1:0}$ for all $N \geq i \geq 1$ using EQ(11.4)
3. Calculating the sums using EQ(11.7)

These steps are illustrated in Figure 11.12. The first and third steps are routine, so most of the attention in the remainder of this section is devoted to alternatives for the group PG logic with different trade-offs between speed, area, and complexity. Some of the hardware can be shared in the bitwise PG logic, as shown in Figure 11.13.

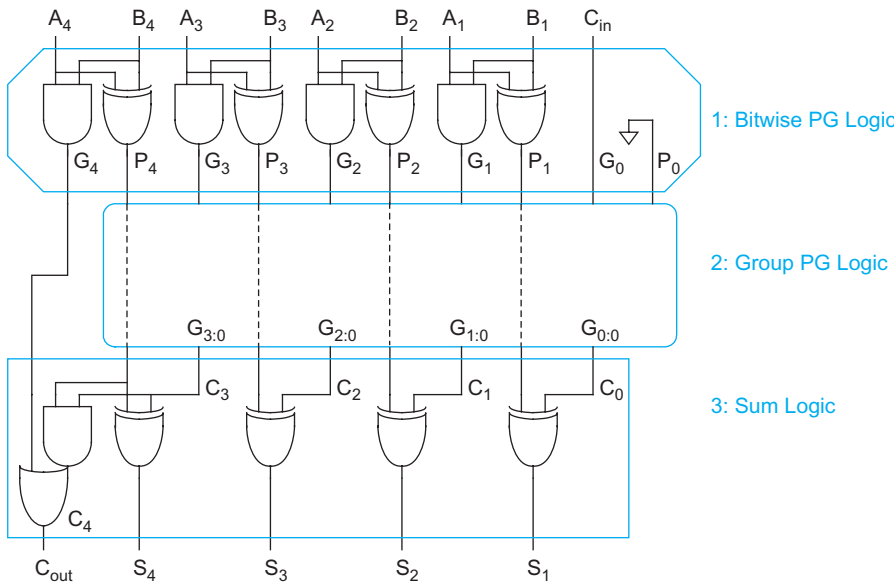


FIGURE 11.12 Addition with generate and propagate logic

Many notations are used in the literature to describe the group PG logic. In general, PG logic is an example of a *prefix* computation [Leighton92]. It accepts inputs $\{P_{N:N}, \dots, P_{0:0}\}$ and $\{G_{N:N}, \dots, G_{0:0}\}$ and computes the *prefixes* $\{G_{N:0}, \dots, G_{0:0}\}$ using the relationship given in EQ(11.4). This relationship is given many names in the literature including the *delta operator*, *fundamental carry operator*, and *prefix operator*. Many other problems such as priority encoding can be posed as prefix computations and all the techniques used to build fast group PG logic will apply, as we will explore in Section 11.10.

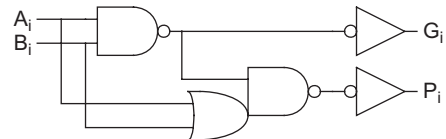


FIGURE 11.13 Shared bitwise PG logic

EQ (11.4) defines *valency-2* (also called *radix-2*) group PG logic because it combines pairs of smaller groups. It is also possible to define higher-valency group logic to use fewer stages of more complex gates [Beaumont-Smith99], as shown in EQ (11.8) and later in Figure 11.16(c). For example, in valency-4 group logic, a group propagates the carry if all four portions propagate. A group generates a carry if the upper portion generates, the second portion generates and the upper propagates, the third generates and the upper two propagate, or the lower generates and the upper three propagate.

$$\left. \begin{aligned} G_{i:j} &= G_{i:k} + P_{i:k} \cdot G_{k-1:l} + P_{i:k} \cdot P_{k-1:l} \cdot G_{l-1:m} + P_{i:k} \cdot P_{k-1:l} \cdot P_{l-1:m} \cdot G_{m-1:j} \\ &= G_{i:k} + P_{i:k} \left(G_{k-1:l} + P_{k-1:l} \left(G_{l-1:m} + P_{l-1:m} G_{m-1:j} \right) \right) \\ P_{i:j} &= P_{i:k} \cdot P_{k-1:l} \cdot P_{l-1:m} \cdot P_{m-1:j} \end{aligned} \right\} \quad (i \geq k > l > m > j) \quad (11.8)$$

Logical Effort teaches us that the best stage effort is about 4. Therefore, it is not necessarily better to build fewer stages of higher-valency gates; simulations or calculations should be done to compare the alternatives for a given process technology and circuit family.

11.2.2.3 PG Carry-Ripple Addition The critical path of the carry-ripple adder passes from carry-in to carry-out along the carry chain majority gates. As the P and G signals will have already stabilized by the time the carry arrives, we can use them to simplify the majority function into an AND-OR gate:³

$$\begin{aligned} C_i &= A_i B_i + (A_i + B_i) C_{i-1} \\ &= A_i B_i + (A_i \oplus B_i) C_{i-1} \\ &= G_i + P_i C_{i-1} \end{aligned} \quad (11.9)$$

Because $C_i = G_{i:0}$, carry-ripple addition can now be viewed as the extreme case of group PG logic in which a 1-bit group is combined with an i -bit group to form an $(i+1)$ -bit group

$$G_{i:0} = G_i + P_i \cdot G_{i-1:0} \quad (11.10)$$

In this extreme, the group propagate signals are never used and need not be computed. Figure 11.14 shows a 4-bit carry-ripple adder. The critical carry path now proceeds through a chain of AND-OR gates rather than a chain of majority gates. Figure 11.15 illustrates the group PG logic for a 16-bit carry-ripple adder, where the AND-OR gates in the group PG network are represented with gray cells.

Diagrams like these will be used to compare a variety of adder architectures in subsequent sections. The diagrams use black cells, gray cells, and white buffers defined in Figure 11.16(a) for valency-2 cells. Black cells contain the group generate and propagate logic (an AND-OR gate and an AND gate) defined in EQ (11.4). Gray cells containing only the group generate logic are used at the final cell position in each column because only the group generate signal is required to compute the sums. Buffers can be used to minimize the load on critical paths. Each line represents a *bundle* of the group generate and propagate signals (propagate signals are omitted after gray cells). The bitwise PG and

³Whenever positive logic such as AND-OR is described, you can also use an AOI gate and alternate positive and negative polarity stages as was done in Figure 11.11(b) to save area and delay.

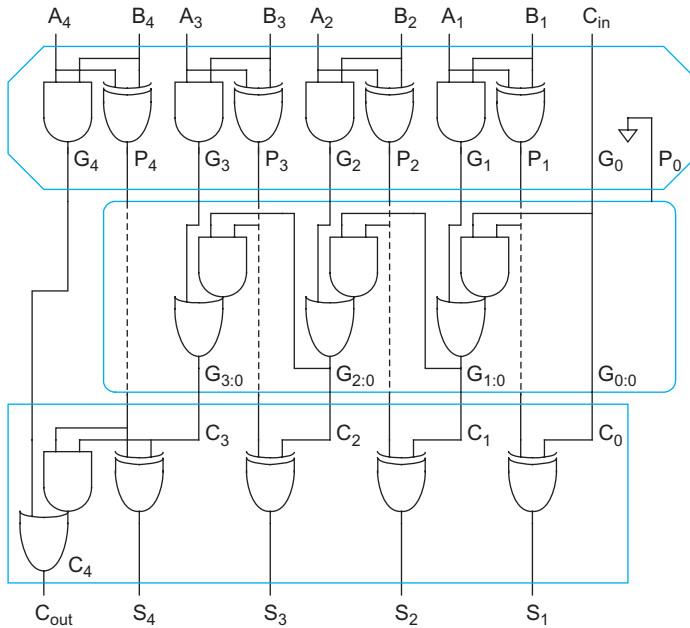


FIGURE 11.14 4-bit carry-ripple adder using PG logic

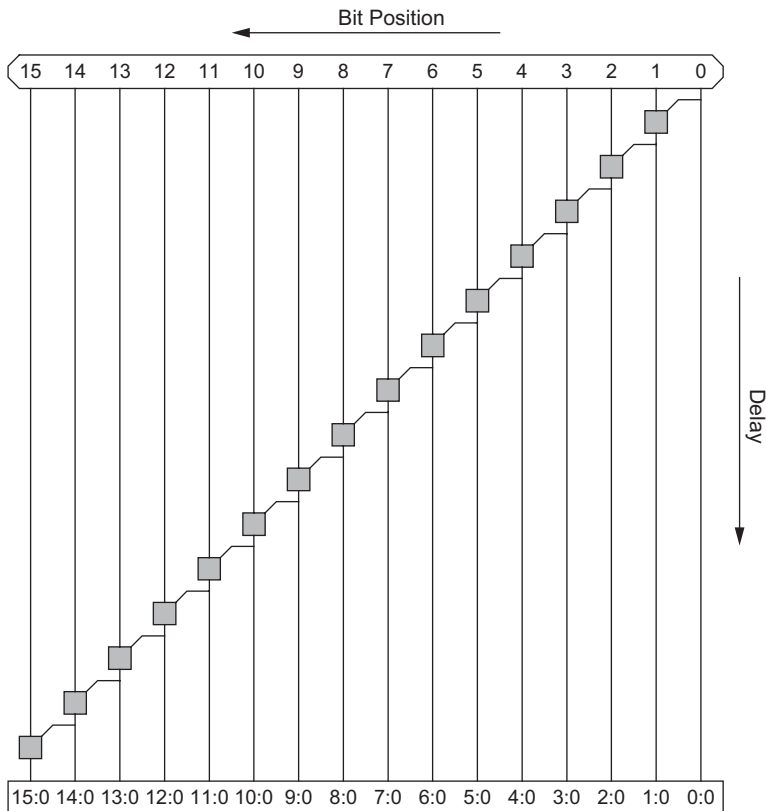


FIGURE 11.15 Carry-ripple adder group PG network

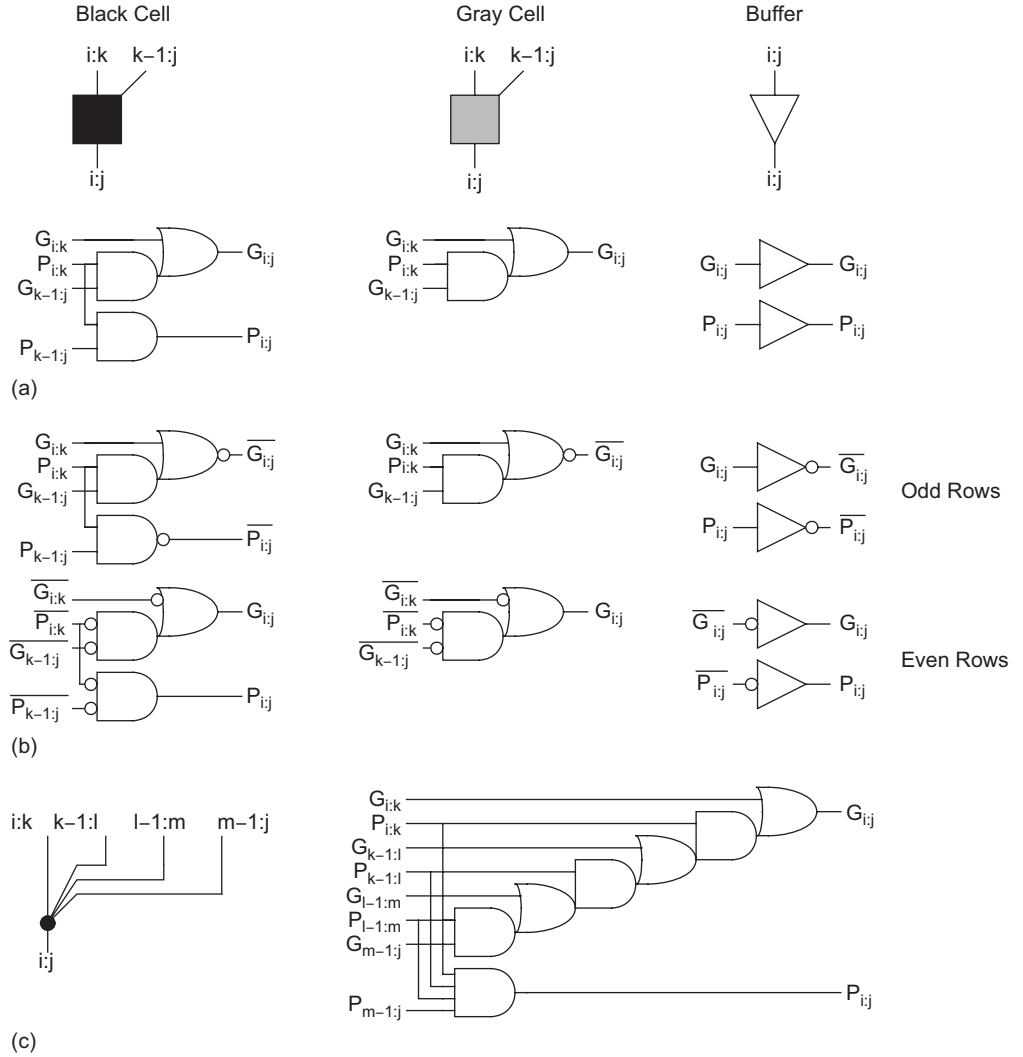


FIGURE 11.16 Group PG cells

sum XORs are abstracted away in the top and bottom boxes and it is assumed that an AND-OR gate operates in parallel with the sum XORs to compute the carry-out:

$$C_{\text{out}} = G_{N:0} = G_N + P_N G_{N-1:0} \quad (11.11)$$

The cells are arranged along the vertical axis according to the time at which they operate [Guyot97]. From Figure 11.15 it can be seen that the carry-ripple adder critical path delay is

$$t_{\text{ripple}} = t_{pg} + (N-1)t_{AO} + t_{\text{xor}} \quad (11.12)$$

where t_{pg} is the delay of the 1-bit propagate/generate gates, t_{AO} is the delay of the AND-OR gate in the gray cell, and t_{xor} is the delay of the final sum XOR. Such a delay estimate is only qualitative because it does not account for fanout or sizing.

Often, using noninverting gates leads to more stages of logic than are necessary. Figure 11.16(b) shows how to alternate two types of inverting stages on alternate rows of the group PG network to remove extraneous inverters. For best performance, G_{k-1j} should drive the inner transistor in the series stack. You can also reduce the number of stages by using higher-valency cells, as shown in Figure 11.16(c) for a valency-4 black cell.

11.2.2.4 Manchester Carry Chain Adder *This section is available in the online Web Enhanced chapter at www.cmosvlsi.com.*



11.2.2.5 Carry-Skip Adder The critical path of CPAs considered so far involves a gate or transistor for each bit of the adder, which can be slow for large adders. The *carry-skip* (also called *carry-bypass*) adder, first proposed by Charles Babbage in the nineteenth century and used for many years in mechanical calculators, shortens the critical path by computing the group propagate signals for each carry chain and using this to skip over long carry ripples [Morgan59, Lehman61]. Figure 11.17 shows a carry skip adder built from 4-bit groups. The rectangles compute the bitwise propagate and generate signals (as in Figure 11.15), and also contain a 4-input AND gate for the propagate signal of the 4-bit group. The skip multiplexer selects the group carry-in if the group propagate is true or the ripple adder carry-out otherwise.

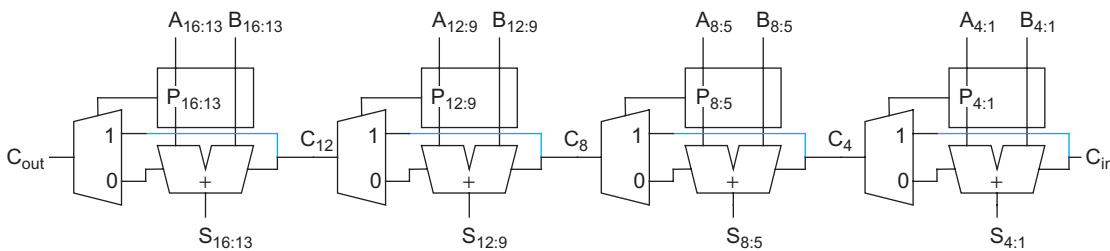


FIGURE 11.17 Carry-skip adder

The critical path through Figure 11.17 begins with generating a carry from bit 1, and then propagating it through the remainder of the adder. The carry must ripple through the next three bits, but then may skip across the next two 4-bit blocks. Finally, it must ripple through the final 4-bit block to produce the sums. This is illustrated in Figure 11.18. The 4-bit ripple chains at the top of the diagram determine if each group generates a carry. The carry skip chain in the middle of the diagram skips across 4-bit blocks. Finally, the 4-bit ripple chains with the blue lines represent *the same adders* that can produce a carry-out when a carry-in is bypassed to them. Note that the final AND-OR and column 16 are not strictly necessary because C_{out} can be computed in parallel with the sum XORs using EQ (11.11).

The critical path of the adder from Figures 11.17 and 11.18 involves the initial PG logic producing a carry out of bit 1, three AND-OR gates rippling it to bit 4, three multiplexers bypassing it to C_{12} , 3 AND-OR gates rippling through bit 15, and a final XOR to produce S_{16} . The multiplexer is an AND22-OR function, so it is slightly slower than the AND-OR function. In general, an N -bit carry-skip adder using k n -bit groups ($N = n \times k$) has a delay of

$$t_{\text{skip}} = t_{pg} + 2(n-1)t_{AO} + (k-1)t_{\text{mux}} + t_{\text{xor}} \quad (11.13)$$

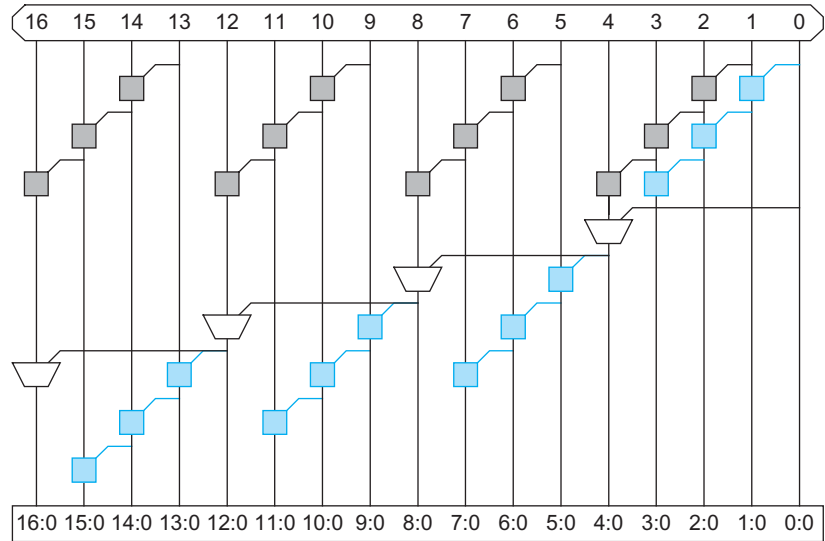


FIGURE 11.18 Carry-skip adder PG network

This critical path depends on the length of the first and last group and the number of groups. In the more significant bits of the network, the ripple results are available early. Thus, the critical path could be shortened by using shorter groups at the beginning and end and longer groups in the middle. Figure 11.19 shows such a PG network using groups of length [2, 3, 4, 4, 3], as opposed to [4, 4, 4, 4], which saves two levels of logic in a 16-bit adder.

The hardware cost of a carry-skip adder is equal to that of a simple carry-ripple adder plus k multiplexers and k n -input AND gates. It is attractive when ripple-carry adders are too slow, but the hardware cost must still be kept low. For long adders, you could use a multilevel skip approach to skip across the skips. A great deal of research has gone into choosing the best group size and number of levels [Majerski67, Oklobdzija85, Guyot87, Chan90, Kantabutra91], although now, parallel prefix adders are generally used for long adders instead.

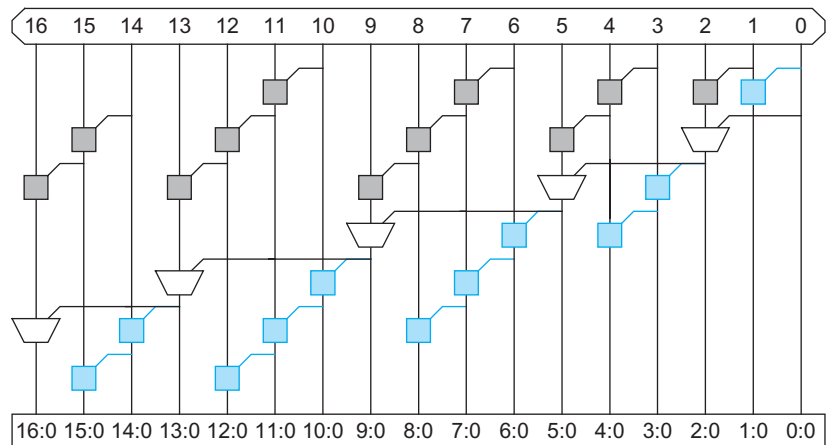


FIGURE 11.19 Variable group size carry-skip adder PG network

It might be tempting to replace each skip multiplexer in Figures 11.17 and 11.18 with an AND-OR gate combining the carry-out of the n -bit adder or the group carry-in and group propagate. Indeed, this works for domino-carry skip adders in which the carry out is precharged each cycle; it also works for carry-lookahead adders and carry-select adders covered in the subsequent section. However, it introduces a sneaky long critical path into an ordinary carry-skip adder. Imagine summing $111\dots111 + 000\dots000 + C_{in}$. All of the group propagate signals are true. If $C_{in} = 1$, every 4-bit block will produce a carry-out. When C_{in} falls, the falling carry signal must ripple through all N bits because of the path through the carry out of each n -bit adder. Domino-carry skip adders avoid this path because all of the carries are forced low during precharge, so they can use AND-OR gates.

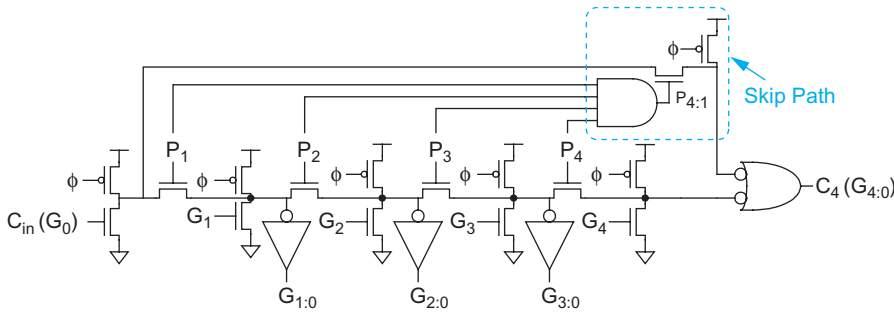


FIGURE 11.20 Carry-skip adder Manchester stage

Figure 11.20 shows how a Manchester carry chain from Section 11.2.2.4 can be modified to perform carry skip [Chan90]. A valency-5 chain is used to skip across groups of 4 bits at a time.

11.2.2.6 Carry-Lookahead Adder The *carry-lookahead adder* (CLA) [Weinberger58] is similar to the carry-skip adder, but computes group generate signals as well as group propagate signals to avoid waiting for a ripple to determine if the first group generates a carry. Such an adder is shown in Figure 11.21 and its PG network is shown in Figure 11.22 using valency-4 black cells to compute 4-bit group PG signals.

In general, a CLA using k groups of n bits each has a delay of

$$t_{cla} = t_{pg} + t_{pg(n)} + [(n-1) + (k-1)]t_{AO} + t_{xor} \quad (11.14)$$

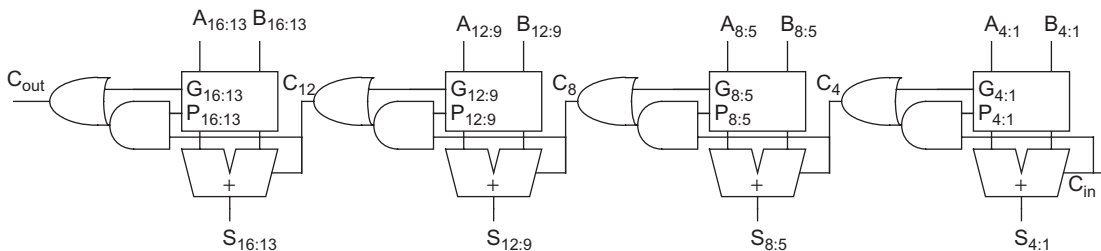


FIGURE 11.21 Carry-lookahead adder

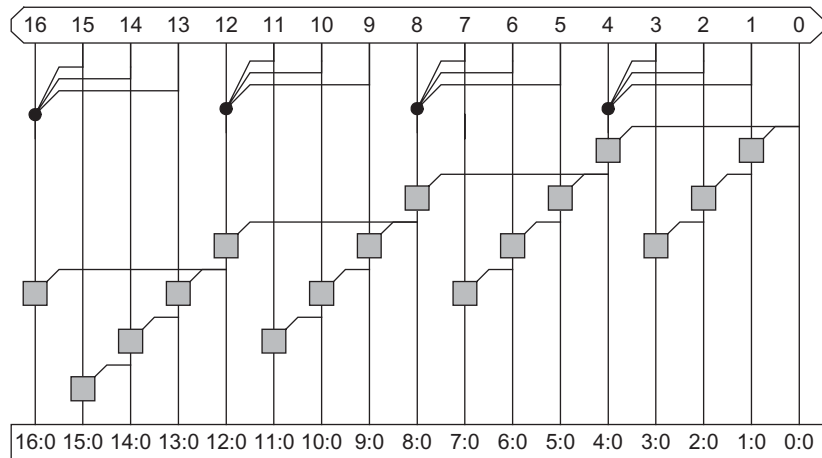


FIGURE 11.22 Carry-lookahead adder group PG network

where $t_{pg(n)}$ is the delay of the AND-OR-AND-OR-...-AND-OR gate computing the valency- n generate signal. This is no better than the variable-length carry-skip adder in Figure 11.19 and requires the extra n -bit generate gate, so the simple CLA is seldom a good design choice. However, it forms the basis for understanding faster adders presented in the subsequent sections.

CLAs often use higher-valency cells to reduce the delay of the n -bit additions by computing the carries in parallel. Figure 11.23 shows such a CLA in which the 4-bit adders are built using Manchester carry chains or multiple static gates operating in parallel.

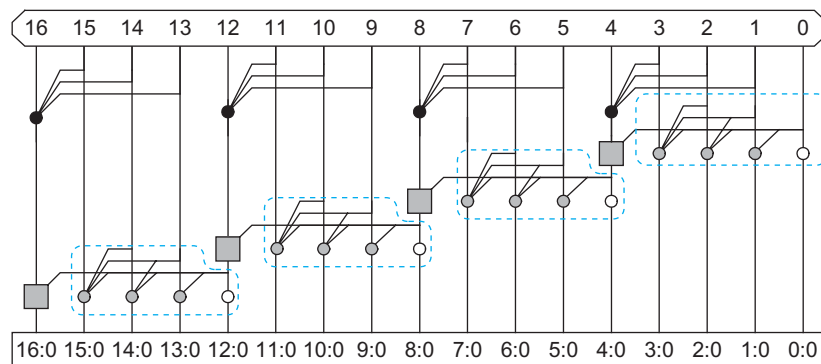


FIGURE 11.23 Improved CLA group PG network

11.2.2.7 Carry-Select, Carry-Increment, and Conditional-Sum Adders The critical path of the carry-skip and carry-lookahead adders involves calculating the carry into each n -bit group, and then calculating the sums for each bit within the group based on the carry-in. A standard logic design technique to accelerate the critical path is to precompute the outputs for both possible inputs, and then use a multiplexer to select between the two output choices. The *carry-select adder* [Bedrij62] shown in Figure 11.24 does this with a pair of

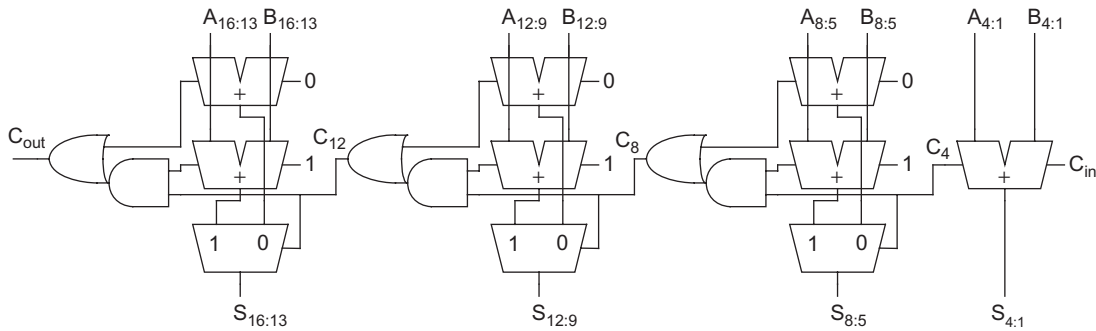


FIGURE 11.24 Carry-select adder

n -bit adders in each group. One adder calculates the sums assuming a carry-in of 0 while the other calculates the sums assuming a carry-in of 1. The actual carry triggers a multiplexer that chooses the appropriate sum. The critical path delay is

$$t_{\text{select}} = t_{pg} + [n + (k - 2)]t_{AO} + t_{\text{mux}} \quad (11.15)$$

The two n -bit adders are redundant in that both contain the initial PG logic and final sum XOR. [Tyagi93] reduces the size by factoring out the common logic and simplifying the multiplexer to a gray cell, as shown in Figure 11.25. This is sometimes called a *carry-increment* adder [Zimmermann96]. It uses a short ripple chain of black cells to compute the PG signals for bits within a group. The bits spanned by each group are annotated on the diagram. When the carry-out from the previous group becomes available, the final gray cells in each column determine the carry-out, which is true if the group generates a carry or if the group propagates a carry and the previous group generated a carry. The carry-increment adder has about twice as many cells in the PG network as a carry-ripple adder. The critical path delay is about the same as that of a carry-select adder because a mux and XOR are comparable, but the area is smaller.

$$t_{\text{increment}} = t_{pg} + \left[(n-1) + (k-1) \right] t_{AO} + t_{\text{xor}} \quad (11.16)$$

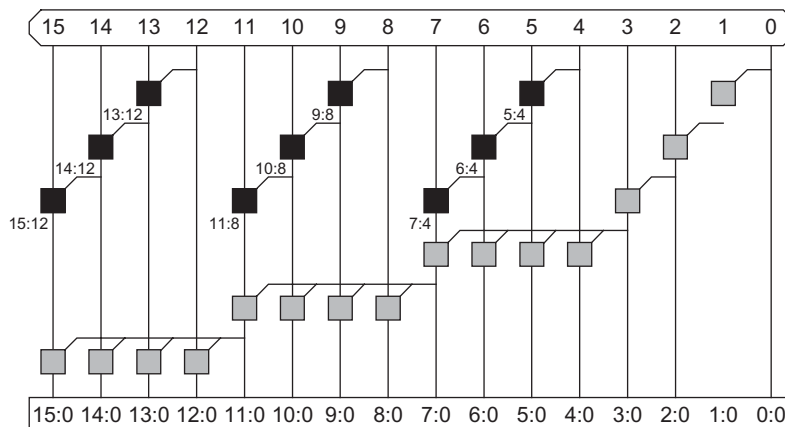


FIGURE 11.25 Carry-increment adder PG network

Of course, Manchester carry chains or higher-valency cells can be used to speed the ripple operation to produce the first group generate signal. In that case, the ripple delay is replaced by a group PG gate delay and the critical path becomes

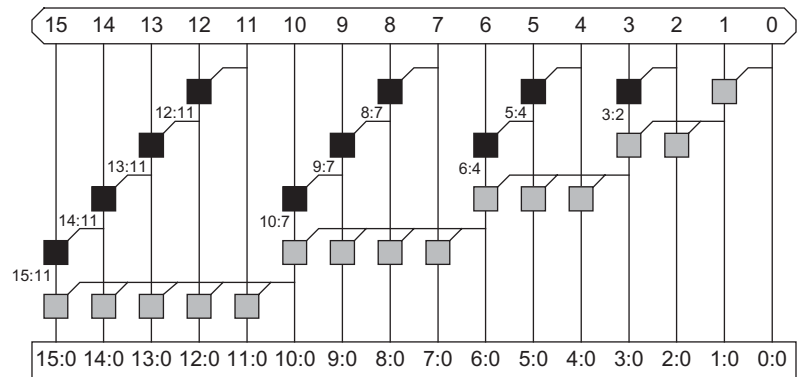
$$t_{\text{increment}} = t_{pg} + t_{pg(n)} + [k-1]t_{AO} + t_{xor} \quad (11.17)$$

As with the carry-skip adder, the carry chains for the more significant bits complete early. Again, we can use variable-length groups to take advantage of the extra time, as shown in Figure 11.26(a). With such a variable group size, the delay reduces to

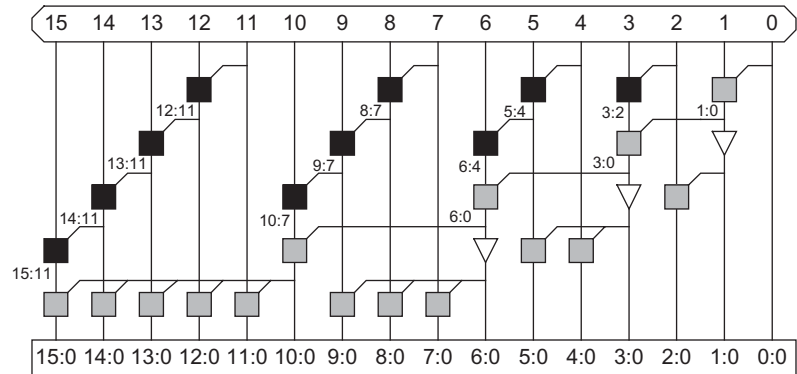
$$t_{\text{increment}} \approx t_{pg} + \sqrt{2N} t_{AO} + t_{xor} \quad (11.18)$$

The delay equations do not account for the fanout that each stage must drive. The fanouts in a variable-length group can become large enough to require buffering between stages. Figure 11.26(b) shows how buffers can be inserted to reduce the branching effort while not impeding the critical lookahead path; this is a useful technique in many other applications.

In wide adders, we can recursively apply multiple levels of carry-select or carry-increment. For example, a 64-bit carry-select or carry-increment. For example, a 64-bit carry-select or carry-increment.



(a)



(b)

FIGURE 11.26 Variable-length carry-increment adder

select adds, each of which selects the carry-in to the next 16-bit group. Taking this to the limit, we obtain the *conditional-sum* adder [Sklansky60] that performs carry-select starting with groups of 1 bit and recursively doubling to $N/2$ bits. Figure 11.27 shows a 16-bit conditional-sum adder. In the first two rows, full adders compute the sum and carry-out for each bit assuming carries-in of 0 and 1, respectively. In the next two rows, multiplexer pairs select the sum and carry-out of the upper bit of each block of two, again assuming carries-in of 0 and 1. In the next two rows, multiplexers select the sum and carry-out of the upper two bits of each block of four, and so forth.

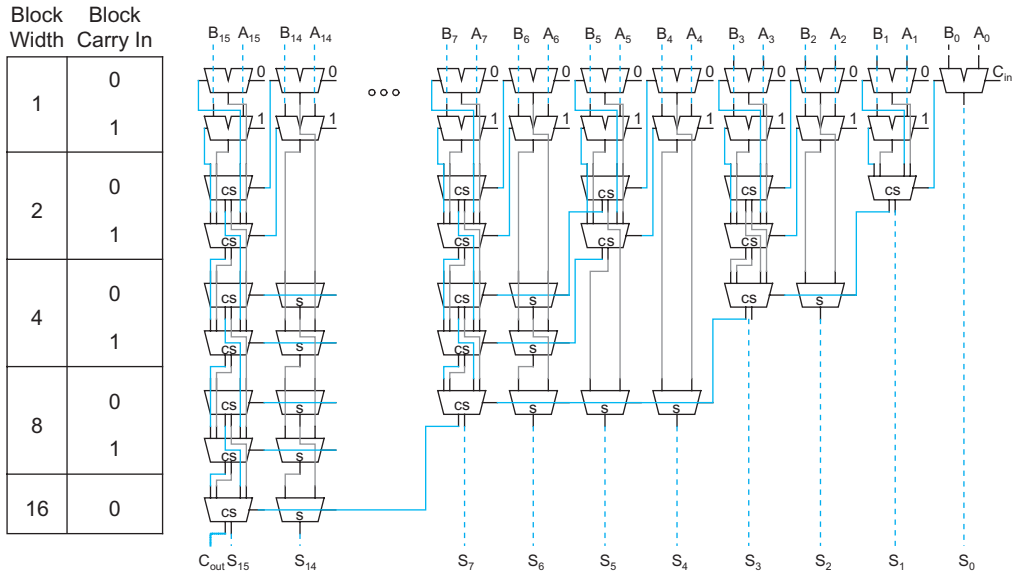


FIGURE 11.27 Conditional-sum adder

Figure 11.28 shows the operation of a conditional-sum adder in action for $N = 16$ with $C_{in} = 0$. In the block width 1 row, a pair of full adders compute the sum and carry-out for each column. One adder operates assuming the carry-in to that column is 0, while the other assumes it is 1. In the block width 2 row, the adder selects the sum for the upper half of each block (the even-numbered columns) based on the carry-out of the lower half. It also computes the carry-out of the pair of bits. Again, this is done twice, for both possibilities of carry-in to the block. In the block width 4 row, the adder again selects the sum for the upper half based on the carry-out of the lower half and finds the carry-out of the entire block. This process is repeated in subsequent rows until the 16-bit sum and the final carry-out are selected.

The conditional-sum adder involves nearly $2N$ full adders and $2N \log_2 N$ multiplexers. As with carry-select, the conditional-sum adder can be improved by factoring out the sum XORs and using AND-OR gates in place of multiplexers. This leads us to the Sklansky tree adder discussed in the next section.

11.2.2.8 Tree Adders For wide adders (roughly, $N > 16$ bits), the delay of carry-lookahead (or carry-skip or carry-select) adders becomes dominated by the delay of passing the carry through the lookahead stages. This delay can be reduced by looking ahead across the look-ahead blocks [Weinberger58]. In general, you can construct a multilevel tree of look-ahead

Block Width	Block Carry In		a	1	0	1	1	1	0	1	1	0	1	1	0	1	1	0	1	0	0	
			b	0	0	0	1	1	0	0	1	1	0	1	1	0	1	1	0	1	0	0
			Block Sum and Carry Out																			
				16	15	14	13	12	11	10	9	8	7	6	5	4	3	2	1	C _{in}		
1	0	s	1	0	1	0	0	0	1	0	1	1	0	1	1	0	1	0	1	1		
		c	0	0	0	1	1	0	0	1	0	0	0	1	0	0	1	0	0	0		
	1	s	0	1	0	1	1	1	0	1	0	0	1	0	1	0	0	1	0			
		c	1	0	1	1	1	0	1	1	1	1	1	1	1	1	1	1	1			
2	0	s	1	0	0	0	0	0	0	0	1	1	0	1	0	0	1	1	1			
		c	0	0	1	0	1	1	1	0	0	0	1	1	1	1	1	0	0			
	1	s	1	1	0	1	0	1	0	1	0	0	1	1	0	0	1	1				
		c	0	0	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1			
4	0	s	1	1	0	0	0	1	0	0	0	0	0	1	0	0	1	1	1			
		c	0	0	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1			
	1	s	1	1	0	1	0	1	0	1	0	0	1	0	1	0	0	1	1			
		c	0	0	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1			
8	0	s	1	1	0	1	0	1	0	0	0	0	1	0	0	0	1	1	1			
		c	0	0	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1			
	1	s	1	1	0	1	0	1	0	1	0	1	0	0	1	0	0	1	1			
		c	0	0	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1			
16	0	s	1	1	0	1	0	1	0	1	0	1	0	0	0	1	0	0	1	1	Sum	
		c	0	0	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	C _{out}	
	1	s																				
		c																				

FIGURE 11.28 Conditional-sum addition example

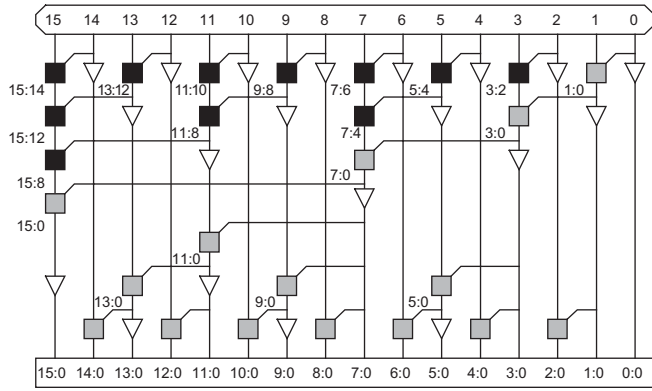
structures to achieve delay that grows with $\log N$. Such adders are variously referred to as *tree* adders, *logarithmic* adders, *multilevel-lookahead* adders, *parallel-prefix* adders, or simply *lookahead* adders. The last name appears occasionally in the literature, but is not recommended because it does not distinguish whether multiple levels of lookahead are used.

There are many ways to build the lookahead tree that offer trade-offs among the number of stages of logic, the number of logic gates, the maximum fanout on each gate, and the amount of wiring between stages. Three fundamental trees are the Brent-Kung, Sklansky, and Kogge-Stone architectures. We begin by examining each in the valency-2 case that combines pairs of groups at each stage.

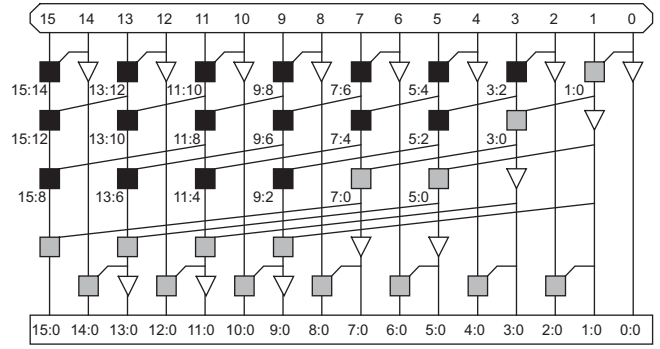
The *Brent-Kung* tree [Brent82] (Figure 11.29(a)) computes prefixes for 2-bit groups. These are used to find prefixes for 4-bit groups, which in turn are used to find prefixes for 8-bit groups, and so forth. The prefixes then fan back down to compute the carries-in to each bit. The tree requires $2\log_2 N - 1$ stages. The fanout is limited to 2 at each stage. The diagram shows buffers used to minimize the fanout and loading on the gates, but in practice, the buffers are generally omitted.

The *Sklansky* or *divide-and-conquer* tree [Sklansky60] (Figure 11.29(b)) reduces the delay to $\log_2 N$ stages by computing intermediate prefixes along with the large group prefixes. This comes at the expense of fanouts that double at each level: The gates fanout to [8, 4, 2, 1] other columns. These high fanouts cause poor performance on wide adders unless the high fanout gates are appropriately sized or the critical signals are buffered before being used for the intermediate prefixes. Transistor sizing can cut into the regularity of the layout because multiple sizes of each cell are required, although the larger gates can spread into adjacent columns. Note that the recursive doubling in the Sklansky tree is analogous to the conditional-sum adder of Figure 11.27. With appropriate buffering, the fanouts can be reduced to [8, 1, 1, 1], as explored in Exercise 11.7.

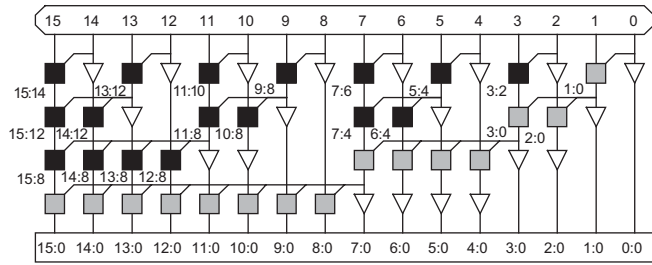
The *Kogge-Stone* tree [Kogge73] (Figure 11.29(c)) achieves both $\log_2 N$ stages and fanout of 2 at each stage. This comes at the cost of many long wires that must be routed between stages. The tree also contains more PG cells; while this may not impact the area if



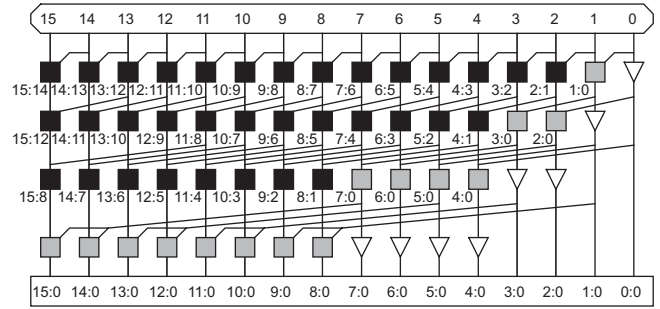
(a) Brent-Kung



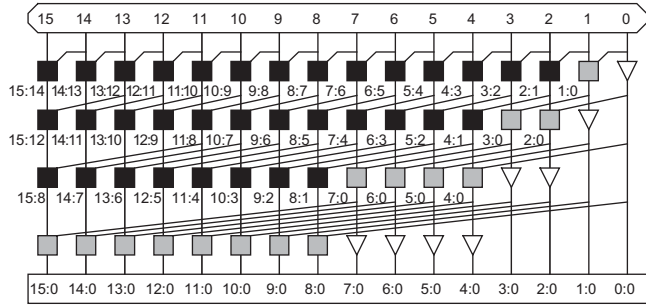
(d) Han-Carlson



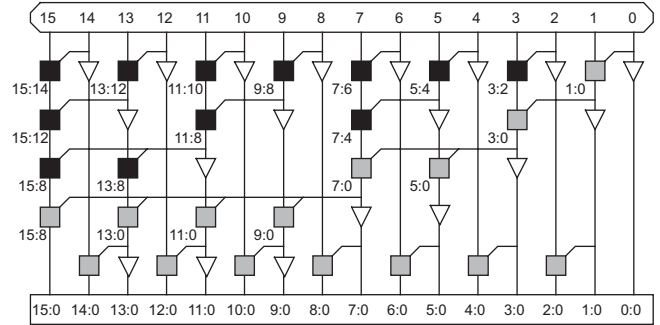
(b) Sklansky



(e) Knowles [2,1,1,1]



(c) Kogge-Stone



(f) Ladner-Fischer

FIGURE 11.29 Tree adder PG networks

the adder layout is on a regular grid, it will increase power consumption. Despite these costs, the Kogge-Stone tree is widely used in high-performance 32-bit and 64-bit adders.

In summary, a Sklansky or Kogge-Stone tree adder reduces the critical path to

$$t_{\text{tree}} \approx t_{pg} + \lceil \log_2 N \rceil t_{AO} + t_{xor} \quad (11.19)$$

An ideal tree adder would have $\log_2 N$ levels of logic, fanout never exceeding 2, and no more than 1 wiring track (G_{ij} and P_{ij} bundle) between each row. The basic tree architectures represent cases that approach the ideal, but each differ in one respect. Brent-Kung

has too many logic levels. Sklansky has too much fanout. And Kogge-Stone has too many wires. Between these three extremes, the Han-Carlson, Ladner-Fischer, and Knowles trees fill out the design space with different compromises between number of stages, fanout, and wire count.

The *Han-Carlson* trees [Han87] are a family of networks between Kogge-Stone and Brent-Kung. Figure 11.29(d) shows such a tree that performs Kogge-Stone on the odd-numbered bits, and then uses one more stage to ripple into the even positions.

The *Knowles* trees [Knowles01] are a family of networks between Kogge-Stone and Sklansky. All of these trees have $\log_2 N$ stages, but differ in the fanout and number of wires. If we say that 16-bit Kogge-Stone and Sklansky adders drive fanouts of $[1, 1, 1, 1]$ and $[8, 4, 2, 1]$ other columns, respectively, the Knowles networks lie between these extremes. For example, Figure 11.29(e) shows a $[2, 1, 1, 1]$ Knowles tree that halves the number of wires in the final track at the expense of doubling the load on those wires.

The *Ladner-Fischer* trees [Ladner80] are a family of networks between Sklansky and Brent-Kung. Figure 11.29(f) is similar to Sklansky, but computes prefixes for the odd-numbered bits and again uses one more stage to ripple into the even positions. Cells at high-fanout nodes must still be sized or ganged appropriately to achieve good speed. Note that some authors use Ladner-Fischer synonymously with Sklansky.

An advantage of the Brent-Kung network and those related to it (Han-Carlson and the Ladner-Fischer network with the extra row) is that for any given row, there is never more than one cell in each pair of columns. These networks have low gate count. Moreover, their layout may be only half as wide, reducing the length of the horizontal wires spanning the adder. This reduces the wire capacitance, which may be a major component of delay in 64-bit and larger adders [Huang00].

Figure 11.30 shows a 3-dimensional taxonomy of the tree adders [Harris03]. If we let $L = \log_2 N$, we can describe each tree with three integers (l, f, t) in the range $[0, L - 1]$. The integers specify the following:

- Logic Levels: $L + l$
- Fanout: $2^f + 1$
- Wiring Tracks: 2^t

The tree adders lie on the plane $l + f + t = L - 1$. 16-bit Brent-Kung, Sklansky, and Kogge-Stone represent vertices of the cube $(3, 0, 0)$, $(0, 3, 0)$ and $(0, 0, 3)$, respectively. Han-Carlson, Ladner-Fischer, and Knowles lie along the diagonals.

11.2.2.9 Higher-Valency Tree Adders Any of the trees described so far can combine more than two groups at each stage [Beaumont-Smith01]. The number of groups combined in each gate is called the *valency* or *radix* of the cell. For example, Figure 11.31 shows 27-bit valency-3 Brent-Kung, Sklansky, Kogge-Stone, and Han-Carlson trees. The rounded boxes mark valency-3 carry chains (that could be constructed using a Manchester carry chain, multiple-output domino gate, or several discrete gates). The trapezoids mark carry-increment operations. The higher-valency designs use fewer stages of logic, but each stage has greater delay. This tends to be a poor trade-off in static CMOS circuits because the stage efforts become much larger than 4, but is good in domino because the logical efforts are much smaller so fewer stages are necessary.

Nodes with large fanouts or long wires can use buffers. The prefix trees can also be internally pipelined for extremely high-throughput operation. Some higher-valency designs combine the initial PG stage with the first level of PG merge. For example, the

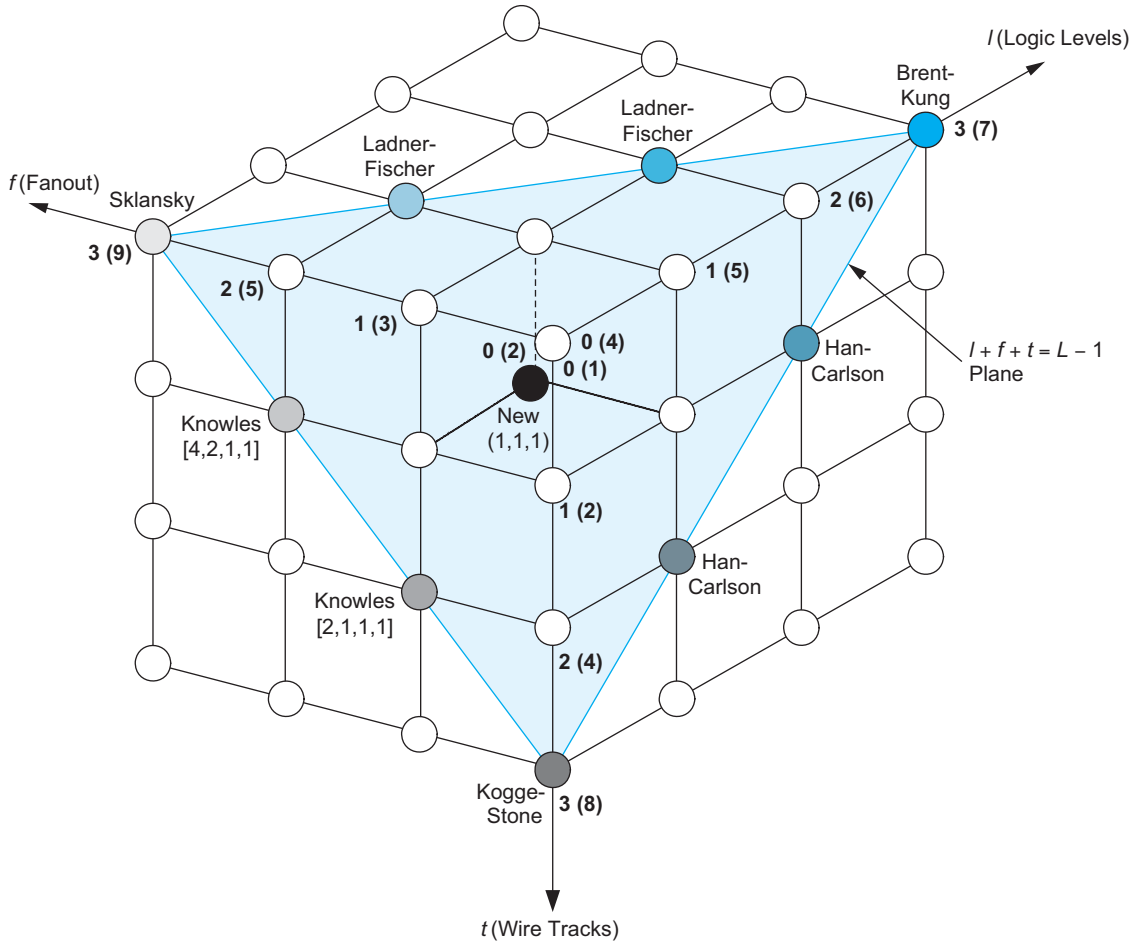


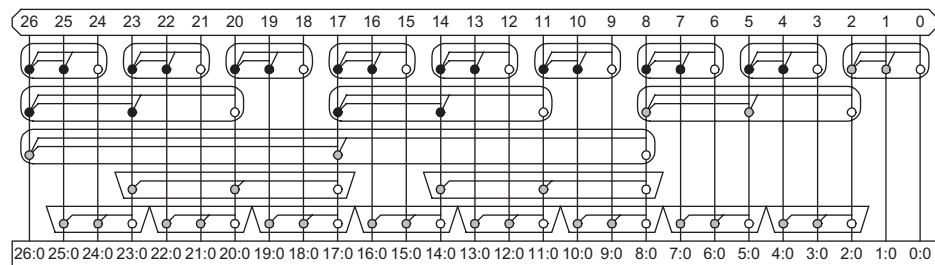
FIGURE 11.30 Taxonomy of prefix networks

Ling adder described in Section 11.2.2.11 computes generate and propagate for up to 4-bit groups from the primary inputs in a single stage.

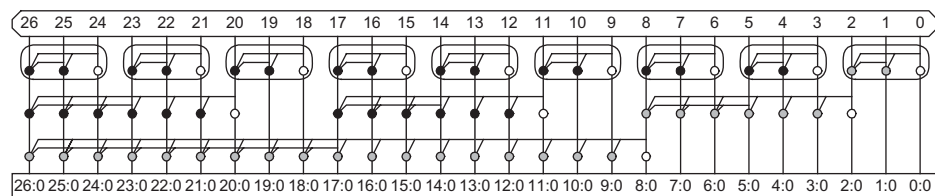
Higher valency (v) adders can still be described in a 3-dimensional taxonomy with $L = \log_v N$ and $l + f + t = L - 1$. There are $L + l$ logic levels, a maximum fanout of $(v - 1)v^f + 1$, and $(v - 1)v^t$ wiring tracks at the worst level.

11.2.2.10 Sparse Tree Adders Building a prefix tree to compute carries in to every bit is expensive in terms of power. An alternative is to only compute carries into short groups (e.g., $s = 2, 4, 8$, or 16 bits). Meanwhile, pairs of s -bit adders precompute the sums assuming both carries-in of 0 and 1 to each group. A multiplexer selects the correct sum for each group based on the carries from the prefix tree. The group length can be balanced such that the carry-in and precomputed sums become available at about the same time. Such a hybrid between a prefix adder and a carry select adder is called a *sparse tree*. s is the *sparse-ness* of the tree.

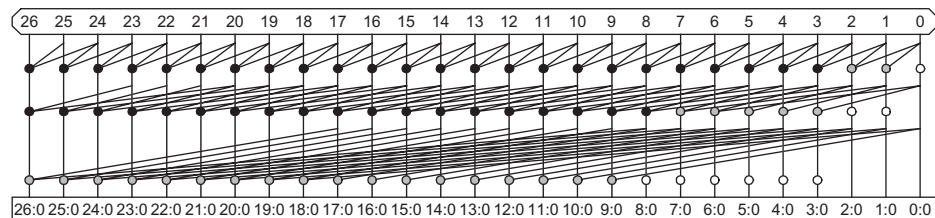
The *spanning-tree adder* [Lynch92] is a sparse tree adder based on a higher-valency Brent-Kung tree of Figure 11.31(a). Figure 11.32 shows a simple valency-3 version that



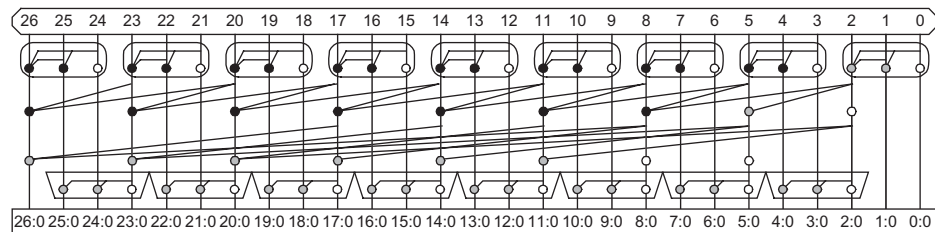
(a) Brent-Kung



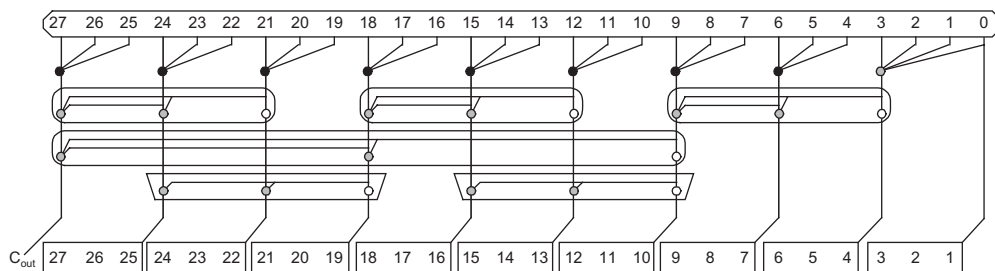
(b) Sklansky



(c) Kogge-Stone



(d) Han-Carlson

FIGURE 11.31 Higher-valency tree adders**FIGURE 11.32** Valency-3 Brent-Kung sparse tree adder with $s = 3$

precomputes sums for $s = 3$ -bit groups and saves one logic level by selecting the output based on the carries into each group. The carry-out (C_{out}) is explicitly shown. Note that the least significant group requires a valency-4 gray cell to compute $G_{3:0}$, the carry-in to the second select block.

[Lynch92] describes a 56-bit spanning-tree design from the AMD AM29050 floating-point unit using valency-4 stages and 8-bit carry select groups. [Kantabutra93] and [Blackburn96] describe optimizing the spanning-tree adder by using variable-length carry-select stages and appropriately selecting transistor sizes.

A carry-select box spanning bits $i \dots j$ is shown in Figure 11.33(a). It uses short carry-ripple adders to precompute the sums assuming carry-in of 0 and 1 to the group, and then selects between them with a multiplexer, as shown in Figure 11.33(b). The adders can be simplified somewhat because the carry-ins are constant, as shown in Figure 11.33(c) for a 4-bit group.

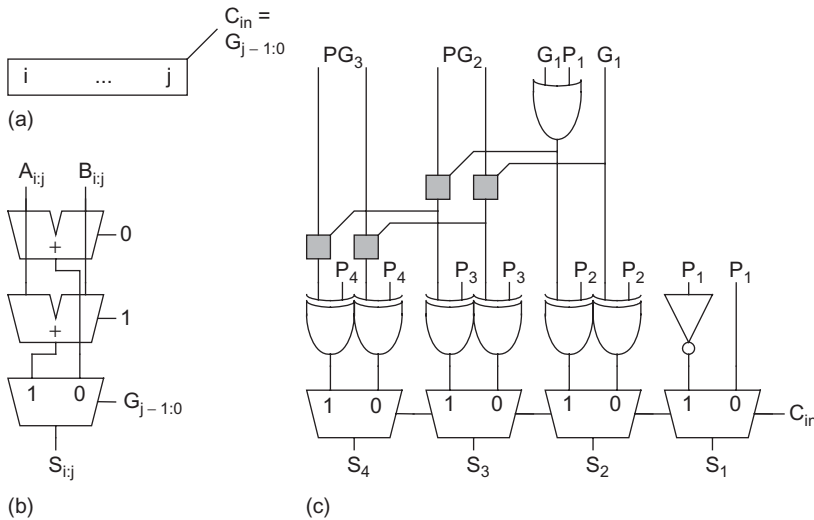
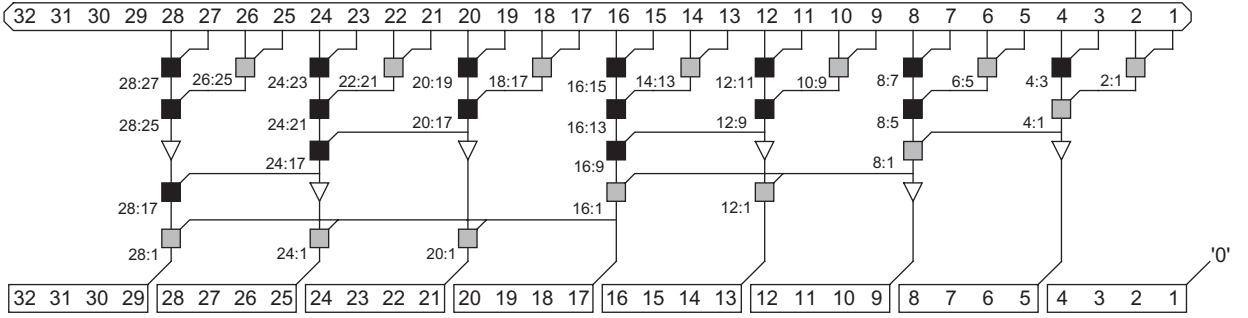
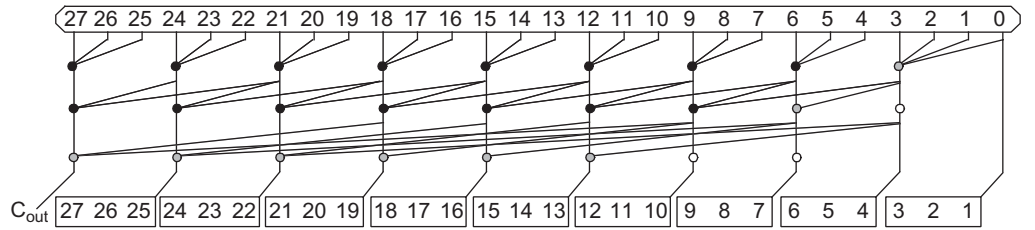


FIGURE 11.33 Carry-select implementation

[Mathew03] describes a 32-bit *sparse-tree adder* using a valency-2 tree similar to Sklansky to compute only the carries into each 4-bit group, as shown in Figure 11.34. This reduces the gate count and power consumption in the tree. The tree can also be viewed as a (2, 2, 0) Ladner-Fischer tree with the final two tree levels and XOR replaced by the select multiplexer. The adder assumes the carry-in is 0 and does not produce a carry-out, saving one input to the least-significant gray box and eliminating the prefix logic in the four most significant columns.

These sparse tree approaches are widely used in high-performance 32–64-bit higher-valency adders because they offer the small number of logic levels of higher-valency trees while reducing the gate count and power consumption in the tree. Figure 11.35 shows a 27-bit valency-3 Kogge-Stone design with carry-select on 3-bit groups. Observe how the number of gates in the tree is reduced threefold. Moreover, because the number of wires is also reduced, the extra area can be used for shielding to reduce path delay. This design can also be viewed as the Han-Carlson adder of Figure 11.31(d) with the last logic level replaced by a carry-select multiplexer.

FIGURE 11.34 Intel valency-2 Sklansky sparse tree adder with $s = 4$ FIGURE 11.35 Valency-3 Kogge-Stone sparse tree adder with $s = 3$

Sparse trees reduce the costly part of the prefix tree. For Kogge-Stone architectures, they reduce the number of wires required by a factor of s . For Sklansky architectures, they reduce the fanout by s . For Brent-Kung architectures, they eliminate the last $\log_v s$ logic levels. In effect, they can move an adder toward the origin in the (l, f, t) design space. These benefits come at the cost of a fanout of s to the final select multiplexer, and of area and power to precompute the sums.



11.2.2.11 Ling Adders Ling discovered a technique to remove one series transistor from the critical group generate path through an adder at the expense of another XOR gate in the sum precomputation [Ling81, Doran88, Bewick94]. The technique depends on using \bar{K} in place of P in the prefix network, and on the observation that $G_i \bar{K}_i = (A_i B_i)(A_i + B_i) = G_i$.

Define a *pseudogenerate* (sometimes called *pseudo-carry*) signal $H_{i:j} = G_i + G_{i-1:j}$. This is simpler than $G_{i:j} = G_i + P_i G_{i-1:j}$. $G_{i:j}$ can be obtained later from $H_{i:j}$ with an AND operation when it is needed:

$$\bar{K}_i H_{i:j} = \bar{K}_i G_i + \bar{K}_i G_{i-1:j} = G_i + \bar{K}_i G_{i-1:j} = G_{i:j} \quad (11.20)$$

The advantage of pseudogenerate signals over regular generate is that the first row in the prefix network is easier to compute.

Also define a *pseudopropagate* signal I that is simply a shifted version of propagate: $I_{i:j} = \bar{K}_{i-1:j-1}$. Group pseudogenerate and pseudopropagate signals are combined using the same black or gray cells as ordinary group generate and propagate signals, as you may show in Exercise 11.11.

$$\begin{aligned} H_{i:j} &= H_{i:k} + I_{i:k} H_{k-1:j} \\ I_{i:j} &= I_{i:k} I_{k-1:j} \end{aligned} \quad (11.21)$$

The true group generate signals are formed from the pseudogenerates using EQ (11.20). These signals can be used to compute the sums with the usual XOR: $S_i = P_i \oplus G_{i-1:0} = P_i \oplus (\bar{K}_{i-1}H_{i-1:0})$. To avoid introducing an AND gate back onto the critical path, we expand S_i in terms of $H_{i-1:0}$

$$S_i = H_{i-1:0} [P_i \oplus \bar{K}_{i-1}] + \bar{H}_{i-1:0} [P_i] \quad (11.22)$$

Thus, sum selection can be performed with a multiplexer choosing either $P_i \oplus \bar{K}_{i-1}$ or P_i based on $H_{i-1:0}$.

The Ling adder technique can be used with any form of adder that uses black and gray cells in a prefix network. It works with any valency and for both domino and static designs. The initial PG stage and the first levels of the prefix network are replaced by a cell that computes the group H and I signals directly. The middle of the prefix network is identical to an ordinary prefix adder but operates on H and I instead of G and P . The sum-selection logic uses the multiplexer from EQ (11.22) rather than an XOR. In sparse trees, the sum out of s -bit blocks is selected directly based on the H signals.

For a valency- v adder, the Ling technique converts a generate gate with v series nMOS transistors and v series pMOS transistors to a pseudogenerate gate with $v - 1$ series nMOS but still v series pMOS. For example, in valency 2, the AOI gate becomes a NOR2 gate. This is not particularly helpful for static logic, but is beneficial for domino implementations because the series pMOS are eliminated and the nMOS stacks are shortened.

Another advantage of the Ling technique is that it allows the first level pseudogenerate and pseudopropagate signals to be computed directly from the A_i and B_i inputs rather than based on G_i and K_i gates. For example, Figure 11.36 compares static gates that compute $G_{2:1}$ and $H_{2:1}$ directly from $A_{2:1}$ and $B_{2:1}$. The H gate has one fewer series transistor and much less parasitic capacitance. $H_{3:1}$ can also be computed directly from $A_{3:1}$ and $B_{3:1}$ using the complex static CMOS gate shown in Figure 11.37(a) [Quach92]. Similarly, Figure 11.37(b) shows a compound domino gate that directly computes $H_{4:1}$ from A and B using only four series transistors rather than the five required for $G_{4:1}$ [Naffziger96, Naffziger98].

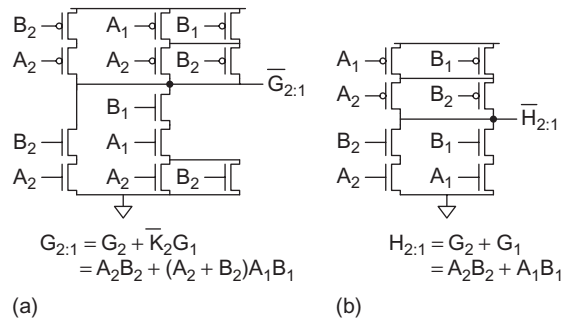


FIGURE 11.36 2-bit generate and pseudogenerate gates using primary inputs

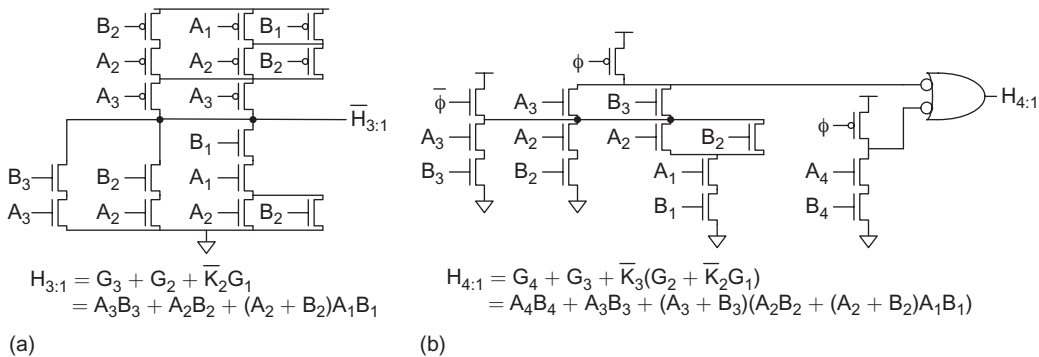


FIGURE 11.37 3-bit and 4-bit pseudogenerate gates using primary inputs

[Jackson04] proposed applying the Ling method recursively to factor out the \bar{K} signal elsewhere in the adder tree. [Burgess09] showed that this recursive Ling technique opens up a new design space containing faster and smaller adders.



11.2.2.12 An Aside on Domino Implementation Issues *This section is available in the online Web Enhanced chapter at www.cmosvlsi.com.*

11.2.2.13 Summary Having examined so many adders, you probably want to know which adder should be used in which application. Table 11.3 compares the various adder architectures that have been illustrated with valency-2 prefix networks. The category “logic levels” gives the number of AND-OR gates in the critical path, excluding the initial PG logic and final XOR. Of course, the delay depends on the fanout and wire loads as well as the number of logic levels. The category “cells” refers to the approximate number of gray and black cells in the network. Carry-lookahead is not shown because it uses higher-valency cells. Carry-select is also not shown because it is larger than carry-increment for the same performance.

In general, carry-ripple adders should be used when they meet timing constraints because they use the least energy and are easy to build. When faster adders are required, carry-increment and carry-skip architectures work well for 8–16 bit lengths. Hybrids combining these techniques are also popular. At word lengths of 32 and especially 64 bits, tree adders are distinctly faster.

TABLE 11.3 Comparison of adder architectures

Architecture	Classification	Logic Levels	Max Fanout	Tracks	Cells
Carry-Ripple		$N - 1$	1	1	N
Carry-Skip ($n = 4$)		$N/4 + 5$	2	1	$1.25N$
Carry-Increment ($n = 4$)		$N/4 + 2$	4	1	$2N$
Carry-Increment (variable group)		$\sqrt{2N}$	$\sqrt{2N}$	1	$2N$
Brent-Kung	$(L-1, 0, 0)$	$2\log_2 N - 1$	2	1	$2N$
Sklansky	$(0, L-1, 0)$	$\log_2 N$	$N/2 + 1$	1	$0.5 N \log_2 N$
Kogge-Stone	$(0, 0, L-1)$	$\log_2 N$	2	$N/2$	$N \log_2 N$
Han-Carlson	$(1, 0, L-2)$	$\log_2 N + 1$	2	$N/4$	$0.5 N \log_2 N$
Ladner Fischer ($l = 1$)	$(1, L-2, 0)$	$\log_2 N + 1$	$N/4 + 1$	1	$0.25 N \log_2 N$
Knowles [2,1,...,1]	$(0, 1, L-2)$	$\log_2 N$	3	$N/4$	$N \log_2 N$

There is still debate about the best tree adder designs; the choice is influenced by power and delay constraints, by domino vs. static and custom vs. synthesis choices, and by the specific manufacturing process. Moreover, careful optimization of a particular architecture is more important than the choice of tree architecture.

When power is no concern, the fastest adders use domino or compound domino circuits [Naffziger96, Park00, Mathew03, Mathew05, Oklobdzija05, Zlatanovici09, Wijeratne07]. Several authors find that the Kogge-Stone architecture gives the lowest

possible delay [Silberman98, Park00, Oklobdzija05, Zlatanovici09]. However, the large number of long wires consume significant energy and require large drivers for speed. Other architectures such as Sklansky [Mathew03] or Han-Carlson [Vangal02] offer better energy efficiency because they have fewer long wires. Valency-4 dynamic gates followed by inverters tend to give a slight speed advantage [Naffziger96, Park00, Zlatanovici09, Harris04, Oklobdzija05], but compound domino implementations using valency-2 dynamic gates followed by valency-2 HI-skew static gates are also used [Mathew03]. Sparse trees save energy in domino adders with little effect on performance [Naffziger96, Mathew03, Zlatanovici09]. The Ling optimization is not used universally, but several studies have found it to be beneficial [Quach92, Naffziger96, Zlatanovici09, Grad04]. The UltraSparc III used a dual-rail domino Kogge-Stone adder [Heald00]. The Itanium 2 and Hewlett Packard PA-RISC lines of 64-bit microprocessors used a dual-rail domino sparse tree Ling adder [Naffziger96, Fetzer02]. The 65 nm Pentium 4 uses a compound domino radix-2 Sklansky sparse tree [Wijeratne07]. A good 64-bit domino adder takes 7–9 FO4 delays and has an area of 4–12 M λ^2 [Naffziger96, Zlatanovici09, Mathew05].

Power-constrained designs use static adders, which consume one third to one tenth the energy of dynamic adders and have a delay of about 13 FO4 [Oklobdzija05, Harris03, Zlatanovici09]. For example, the CELL processor floating point unit uses a valency-2 static Kogge-Stone adder [Oh06].

[Patil07] presents a comprehensive study of energy-delay design space for adders. The paper concludes that the Sklansky architecture is most energy efficient for any delay requirement because it avoids the large number of power-hungry wires in Kogge-Stone and the excessive number of logic levels in Brent-Kung. The high-fanout gates in the Sklansky tree are upsized to maintain a reasonable logical effort. Static adders are most efficient using valency-2 cells, which provide a stage effort of about 4. Domino adders are most efficient alternating valency-4 dynamic gates with static inverters. The sum precomputation logic in a static sparse tree adder costs more energy than it saves from the prefix network. In a domino adder, a sparseness of 2 does save energy because the sum precomputation can be performed with static gates. Figure 11.38 shows some results, finding that static adders are most energy-efficient for slow adders, while domino become better at high speed requirements and dual-rail domino Ling adders are preferable only for the very fastest and most energy-hungry adders. The very fast delays are achieved using a higher V_{DD} and lower V_t . [Zlatanovici09] explores the energy-delay space for 64-bit domino adders and came to the contradictory conclusion that Kogge-Stone is superior. Again, alternating valency-4 dynamic gates with static inverters and using a sparseness of 2 gave the best results, as shown in Figure 11.39. Other reasonable adders are almost as good in the energy-delay space, so there is not a compelling reason to choose one topology over another and the debate about the “best” adder will doubtlessly rage into the future.

Good logic synthesis tools automatically map the “+” operator onto an appropriate adder to meet timing constraints while minimizing area. For example, the Synopsys DesignWare libraries contain carry-ripple adders, carry-select adders, carry-lookahead adders, and a variety of prefix adders. Figure 11.40 shows the results of synthesizing

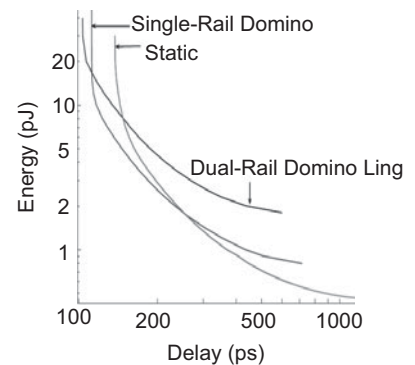


FIGURE 11.38 Energy-delay trade-offs for 90 nm 32-bit Sklansky static, domino, and dual-rail domino adders. FO4 inverter delay in this process at 1.0 V and nominal V_t is 31 ps. (© IEEE 2007.)

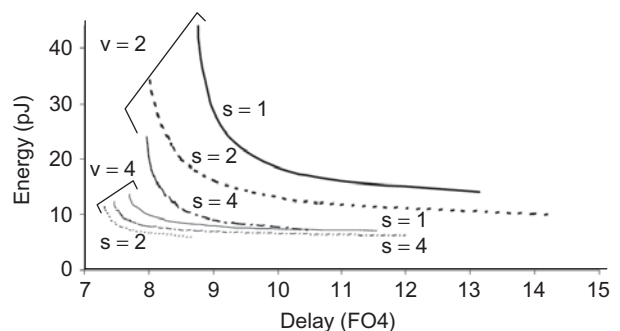


FIGURE 11.39 Energy-delay trade-offs for 90 nm 64-bit domino Kogge-Stone Ling adders as a function of valency (v) and sparseness (s). (© IEEE 2009.)

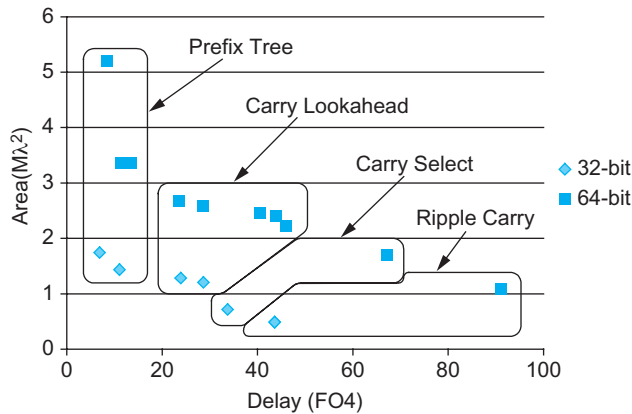


FIGURE 11.40 Area vs. delay of synthesized adders

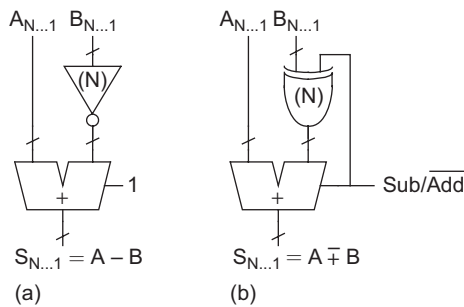


FIGURE 11.41 Subtractors

11.2.3 Subtraction

An N -bit subtracter uses the two's complement relationship

$$A - B = A + \bar{B} + 1 \quad (11.23)$$

This involves inverting one operand to an N -bit CPA and adding 1 via the carry input, as shown in Figure 11.41(a). An adder/subtractor uses XOR gates to conditionally invert B , as shown in Figure 11.41(b). In prefix adders, the XOR gates on the B inputs are sometimes merged into the bitwise PG circuitry.

11.2.4 Multiple-Input Addition

The most obvious method of adding k N -bit words is with $k - 1$ cascaded CPAs as illustrated in Figure 11.42(a) for $0001 + 0111 + 1101 + 0010$. This approach consumes a large amount of hardware and is slow. A better technique is to note that a full adder sums three inputs of unit weight and produces a sum output of unit weight and a carry output of double weight. If N full adders are used in parallel, they can accept three N -bit input words $X_{N...1}$, $Y_{N...1}$, and $Z_{N...1}$, and produce two N -bit output words $S_{N...1}$ and $C_{N...1}$, satisfying $X + Y + Z = S + 2C$, as shown in Figure 11.42(b). The results correspond to the sums and carries-out of each adder. This is called *carry-save redundant format* because the carry outputs are preserved rather than propagated along the adder. The full adders in this application are sometimes called *[3:2] carry-save adder (CSA)* because they accept three inputs and produce two outputs in carry-save form. When the carry word C is shifted left by one position (because it has double weight) and added to the sum word S with an ordinary CPA, the result is $X + Y + Z$. Alternatively, a fourth input word can be added to the carry-save redundant result with another row of CSAs, again resulting in a carry-save redundant result. Such carry-save addition of four numbers is illustrated in Figure 11.42(c), where the underscores in the carry outputs serve as reminders that the carries must be shifted left one column on account of their greater weight.

The critical path through a [3:2] adder is for the sum computation, which involves one 3-input XOR, or two levels of XOR2. This is much faster than a CPA. In general, k

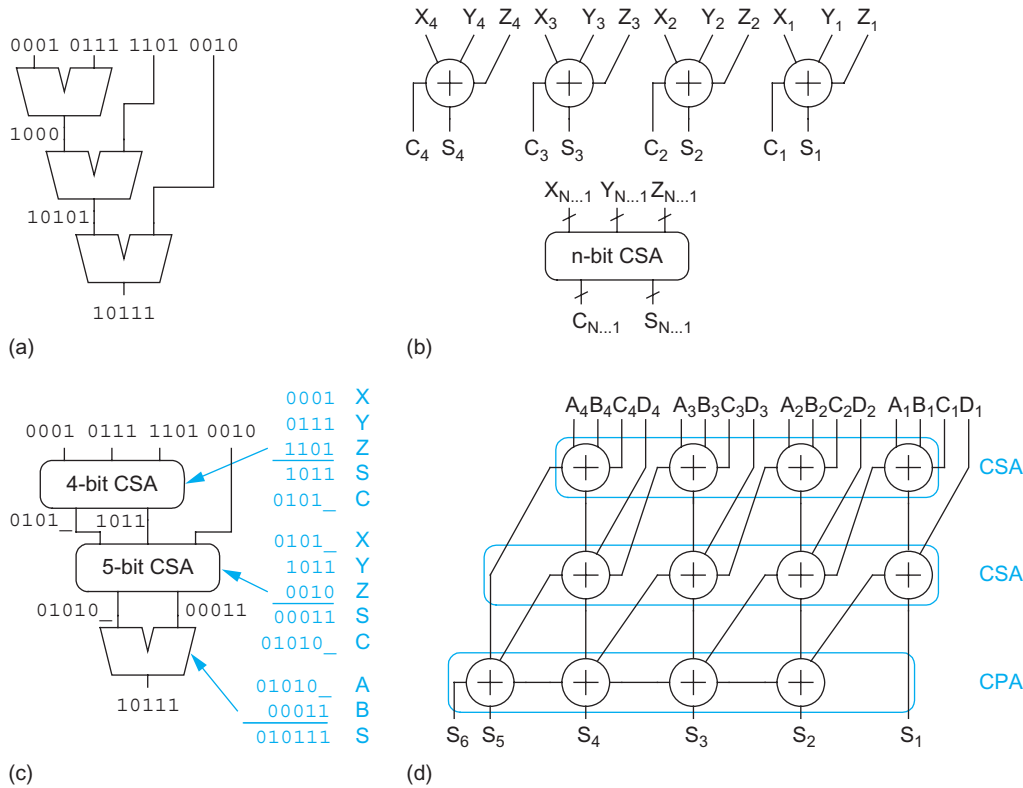


FIGURE 11.42 Multiple-input adders

numbers can be summed with $k - 2$ [3:2] CSAs and only one CPA. This approach will be exploited in Section 11.9 to add many partial products in a multiplier rapidly. The technique dates back to von Neumann's early computer [Burks46].

When one of the inputs to a CSA is a constant, the hardware can be reduced further. If a bit of the input is 0, the CSA column reduces to a half-adder. If the bit is 1, the CSA column simplifies to $S = A \oplus B$ and $C = A + B$.

11.2.5 Flagged Prefix Adders

Sometimes it is necessary to compute either $A + B$, and then, depending on a late-arriving control signal, adding 1. Some applications include calculating $A + B \bmod 2^n - 1$ for cryptography and Reed-Solomon coding, computing the absolute difference $|A - B|$, doing addition/subtraction of sign-magnitude numbers, and performing rounding in certain floating-point adders [Beaumont-Smith99]. A straightforward approach is to build two adders, provide a carry to one, and select between the results. [Burgess02] describes a clever alternative called a *flagged prefix adder* that uses much less hardware.

A flagged prefix adder receives A , B , and a control signal, inc , and computes $A + B + inc$. Recall that an ordinary adder computes the prefixes $G_{i-1:0}$ as the carries into each column i , then computes the sum $S_i = P_i \oplus G_{i-1:0}$. In this situation, there is no C_{in} and hence column 0 is omitted; $G_{i-1:1}$ is used instead. The goal of the flagged prefix adder is to adjust these carries when inc is asserted. A flagged prefix adder instead uses



$G'_{i-1:1} = G_{i-1:1} + P_{i-1:1} \cdot inc$. Thus, if inc is true, it generates a carry into all of the low order bits whose group propagate signals are TRUE. The modified prefixes, $G'_{i-1:1}$, are called *flags*. The sums are computed in the same way with an XOR gate: $S_i = P_i \oplus G'_{i-1:1}$.

To produce these flags, the flagged prefix adder uses one more row of gray cells. This requires that the former bottom row of gray cells be converted to black cells to produce the group propagate signals. Figure 11.43 shows a flagged prefix Kogge–Stone adder. The new row, shown in blue, is appended to perform the late increment. Column 0 is eliminated because there is no C_{in} , but column 16 is provided because applications of flagged adders will need the generate and propagate signals spanning the entire n bits.

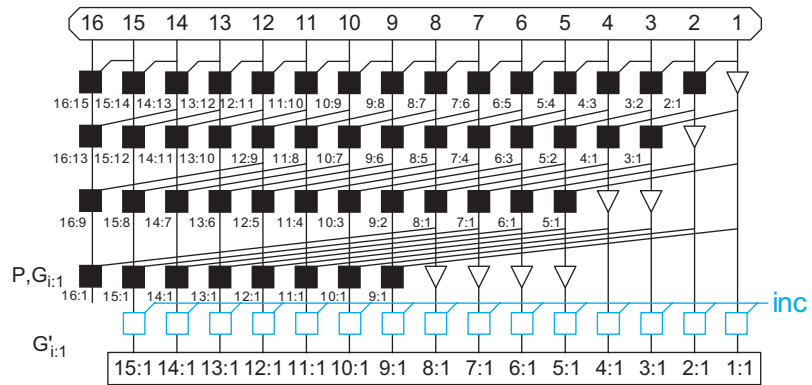


FIGURE 11.43 Flagged prefix Kogge–Stone adder

11.2.5.1 Modulo $2^n - 1$ Addition To compute $A + B \bmod 2^n - 1$ for unsigned operands, an adder should first compute $A + B$. If the sum is greater than or equal to $2^n - 1$, the result should be incremented and truncated back to n bits. $G_{n:1}$ is TRUE if the adder will overflow; i.e., the result is greater than $2^n - 1$. $P_{n:1}$ is TRUE if all columns propagate, which only occurs when the sum equals $2^n - 1$. Hence, modular addition can be done with a flagged prefix adder using $inc = G_{n:1} + P_{n:1}$.

Compared to ordinary addition, modular addition requires one more row of black cells, an OR gate to compute inc , and a buffer to drive inc across all n bits.

11.2.5.2 Absolute Difference $|A - B|$ is called the *absolute difference* and is commonly used in applications such as video compression. The most straightforward approach is to compute both $A - B$ and $B - A$, then select the positive result. A more efficient technique is to compute $A + \bar{B}$ and look at the sign, indicated by $\bar{G}_{n:1}$. If the result is negative, it should be inverted to obtain $B - A$. If the result is positive, it should be incremented to obtain $A - B$.

All of these operations can be performed using a flagged prefix adder enhanced to invert the result conditionally. Modify the sum logic to calculate $S_i = (P_i \oplus inv) \oplus G'_{i-1:1}$. Choose $inv = \bar{G}_{n:1}$ and $inc = G_{n:1}$.

Compared to ordinary addition, absolute difference requires a bank of inverters to obtain \bar{B} , one more row of black cells, buffers to drive inv and inc across all n bits, and a row of XORs to invert the result conditionally. Note that $(P_i \oplus inv)$ can be precomputed so this does not affect the critical path.

11.2.5.3 Sign-Magnitude Arithmetic Addition of sign-magnitude numbers involves examining the signs of the operands. If the signs agree, the magnitudes are added and the

sign is unchanged. If the signs differ, the absolute difference of the magnitudes must be computed. This can be done using the flagged carry adder described in the previous section. The sign of the result is $\text{sign}(A) \oplus G_{n:1}$.

Subtraction is identical except that the sign of B is first flipped.

11.3 One/Zero Detectors

Detecting all ones or zeros on wide N -bit words requires large fan-in AND or NOR gates. Recall that by DeMorgan's law, AND, OR, NAND, and NOR are fundamentally the same operation except for possible inversions of the inputs and/or outputs. You can build a tree of AND gates, as shown in Figure 11.44(a). Here, alternate NAND and NOR gates have been used. The path has $\log N$ stages. In general, the minimum logical effort is achieved with a tree alternating NAND gates and inverters and the path logical effort is

$$G_{\text{and}}(N) = \left(\frac{4}{3}\right)^{\log_2 N} = N^{\log_2 \frac{4}{3}} = N^{0.415} \quad (11.24)$$

A rough estimate of the path delay driving a path electrical effort of H using static CMOS gates is

$$D \approx (\log_4 F) t_{FO4} = (\log_4 H + 0.415 \log_4 N) t_{FO4} \quad (11.25)$$

where t_{FO4} is the fanout-of-4 inverter delay.

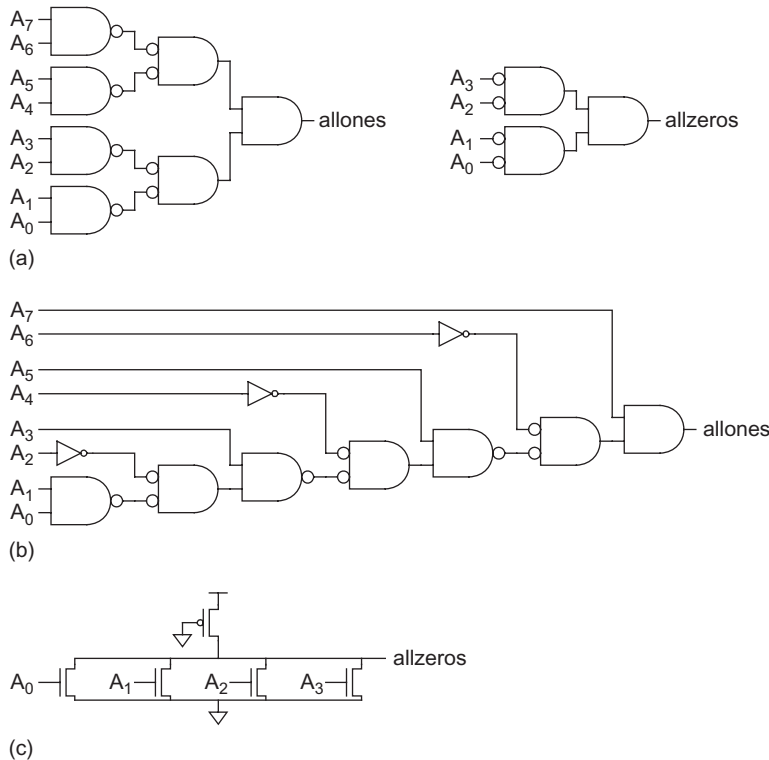


FIGURE 11.44 One/zero detectors

If the word being checked has a natural skew in the arrival time of the bits (such as at the output of a ripple adder), the designer might consider an asymmetric design that favors the late-arriving inputs, as shown in Figure 11.44(b). Here, the delay from the latest bit A_7 is a single gate.

Another fast detector uses a pseudo-nMOS or dynamic NOR structure to perform the “wired-OR,” as shown in Figure 11.44(c). This works well for words up to about 16 bits; for larger words, the gates can be split into 8–16-bit chunks to reduce the parasitic delay and avoid problems with subthreshold leakage.

11.4 Comparators

11.4.1 Magnitude Comparator

A *magnitude comparator* determines the larger of two binary numbers. To compare two unsigned numbers A and B , compute $B - A = B + \bar{A} + 1$. If there is a carry-out, $A \leq B$; otherwise, $A > B$. A zero detector indicates that the numbers are equal. Figure 11.45 shows a 4-bit unsigned comparator built from a carry-ripple adder and two’s complementer. The relative magnitude is determined from the carry-out (C) and zero (Z) signals according to Table 11.4. For wider inputs, any of the faster adder architectures can be used.

Comparing signed two’s complement numbers is slightly more complicated because of the possibility of overflow when subtracting two numbers with different signs. Instead of simply examining the carry-out, we must determine if the result is negative (N , indicated by the most significant bit of the result) and if it overflows the range of possible signed numbers. The overflow signal V is true if the inputs had different signs (most significant bits) and the output sign is different from the sign of B . The actual sign of the difference $B - A$ is $S = N \oplus V$ because overflow flips the sign. If this corrected sign is negative ($S = 1$), we know $A > B$. Again, the other relations can be derived from the corrected sign and the Z signal.

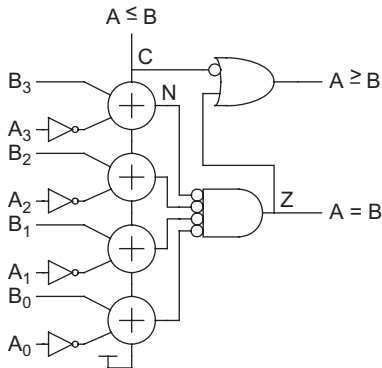


FIGURE 11.45

Unsigned magnitude comparator

TABLE 11.4 Magnitude comparison

Relation	Unsigned Comparison	Signed Comparison
$A = B$	Z	Z
$A \neq B$	\bar{Z}	\bar{Z}
$A < B$	$C \cdot \bar{Z}$	$\bar{S} \cdot \bar{Z}$
$A > B$	C	S
$A \leq B$	C	\bar{S}
$A \geq B$	$\bar{C} + Z$	$S + Z$

11.4.2 Equality Comparator

An *equality comparator* determines if ($A = B$). This can be done more simply and rapidly with XNOR gates and a ones detector, as shown in Figure 11.46.

11.4.3 $K = A + B$ Comparator

Sometimes it is necessary to determine if $(A + B = K)$. For example, the sum-addressed memory [Heald98] described in Section 12.2.2.4 contains a decoder that must match against the sum of two numbers, such as a register base address and an immediate offset. Remarkably, this comparison can be done faster than computing $A + B$ because no carry propagation is necessary. The key is that if you know A and B , you also know what the carry into each bit must be if $K = A + B$ [Cortadella92]. Therefore, you only need to check adjacent pairs of bits to verify that the previous bit produces the carry required by the current bit, and then use a ones detector to check that the condition is true for all N pairs. Specifically, if $K = A + B$, Table 11.5 lists what the carry-in c_{i-1} must have been for this to be true and what the carry-out c_i will be for each bit position i .

TABLE 11.5 Required and generated carries if $K = A + B$

A_i	B_i	K_i	c_{i-1} (required)	c_i (produced)
0	0	0	0	0
0	0	1	1	0
0	1	0	1	1
0	1	1	0	0
1	0	0	1	1
1	0	1	0	0
1	1	0	0	1
1	1	1	1	1

From this table, you can see that the required c_{i-1} for bit i is

$$c_{i-1} = A_i \oplus B_i \oplus K_i \quad (11.26)$$

and the c_{i-1} produced by bit $i-1$ is

$$c_{i-1} = (A_{i-1} \oplus B_{i-1}) \bar{K}_{i-1} + A_{i-1} \cdot B_{i-1} \quad (11.27)$$

Figure 11.47 shows one bitslice of a circuit to perform this operation. The XNOR gate is used to make sure that the required carry matches the produced carry at each bit position; then the AND gate checks that the condition is satisfied for all bits.

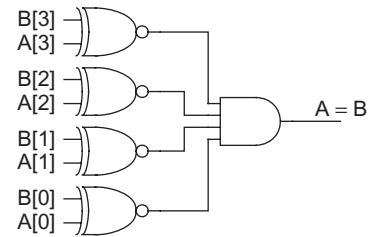


FIGURE 11.46 Equality comparator

11.5 Counters

Two commonly used types of counters are *binary counters* and *linear-feedback shift registers*. An N -bit binary counter sequences through 2^N outputs in binary order. Simple designs have a minimum cycle time that increases with N , but faster designs operate in constant time. An N -bit linear-feedback shift register sequences through up to $2^N - 1$ outputs in pseudo-random order. It has a short minimum cycle time independent of N , so it is useful for extremely fast counters as well as pseudo-random number generation.

asserted for a cycle to load the new *count* from the *shadow* register rather than the adder (which may not have had enough time to ripple carries).



11.5.3 Ring and Johnson Counters

A *ring counter* consists of an M -bit shift register with the output fed back to the input, as shown in Figure 11.53(a). On reset, the first bit is initialized to 1 and the others are initialized to 0. *TC* toggles once every M cycles. Ring counters are a convenient way to build extremely fast prescalars because there is no logic between flip-flops, but they become costly for larger M .

A *Johnson* or *Mobius counter* is similar to a ring counter, but inverts the output before it is fed back to the input, as shown in Figure 11.53(b). The flip-flops are reset to all zeros and count through $2M$ states before repeating. Table 11.6 shows the sequence for a 3-bit Johnson counter.

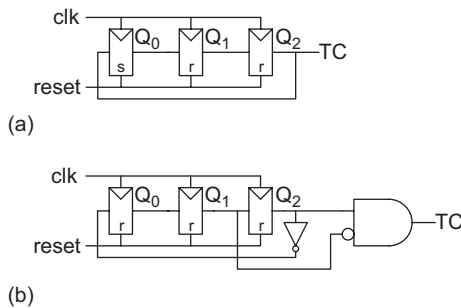


FIGURE 11.53 3-bit ring and Johnson counters

TABLE 11.6 Johnson counter sequence

Cycle	Q_0	Q_1	Q_2	TC
0	0	0	0	0
1	1	0	0	0
2	1	1	0	0
3	1	1	1	0
4	0	1	1	0
5	0	0	1	1
6	0	0	0	0
Repeats forever				

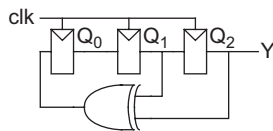


FIGURE 11.54 3-bit LFSR

11.5.4 Linear-Feedback Shift Registers

A linear-feedback shift register (LFSR) consists of N registers configured as a shift register. The input to the shift register comes from the XOR of particular bits of the register, as shown in Figure 11.54 for a 3-bit LFSR. On reset, the registers must be initialized to a nonzero value (e.g., all 1s). The pattern of outputs for the LFSR is shown in Table 11.7.

TABLE 11.7 LFSR sequence

Cycle	Q_0	Q_1	Q_2 / Y
0	1	1	1
1	0	1	1
2	0	0	1
3	1	0	0
4	0	1	0
5	1	0	1
6	1	1	0
7	1	1	1
Repeats forever			

This LFSR is an example of a *maximal-length* shift register because its output sequences through all $2^n - 1$ combinations (excluding all 0s). The inputs fed to the XOR are called the *tap sequence* and are often specified with a *characteristic polynomial*. For example, this 3-bit LFSR has the characteristic polynomial $1 + x^2 + x^3$ because the taps come after the second and third registers.

The output Y follows the 7-bit sequence [1110010]. This is an example of a *pseudo-random bit sequence* (PRBS). LFSRs are used for high-speed counters and pseudo-random number generators. The pseudo-random sequences are handy for built-in self-test and bit-error-rate testing in communications links. They are also used in many spread-spectrum communications systems such as GPS and CDMA where their correlation properties make other users look like uncorrelated noise.

Table 11.8 lists characteristic polynomials for some commonly used maximal-length LFSRs. For certain lengths, N , more than two taps may be required. For many values of N , there are multiple polynomials resulting in different maximal-length LFSRs. Observe that the cycle time is set by the register and a small number of XOR delays. [Golomb81] offers the definitive treatment on linear-feedback shift registers.

TABLE 11.8 Characteristic polynomials

N	Polynomial
3	$1 + x^2 + x^3$
4	$1 + x^3 + x^4$
5	$1 + x^3 + x^5$
6	$1 + x^5 + x^6$
7	$1 + x^6 + x^7$
8	$1 + x^1 + x^6 + x^7 + x^8$
9	$1 + x^5 + x^9$
15	$1 + x^{14} + x^{15}$
16	$1 + x^4 + x^{13} + x^{15} + x^{16}$
23	$1 + x^{18} + x^{23}$
24	$1 + x^{17} + x^{22} + x^{23} + x^{24}$
31	$1 + x^{28} + x^{31}$
32	$1 + x^{10} + x^{30} + x^{31} + x^{32}$

Example 11.1

Sketch an 8-bit linear-feedback shift register. How long is the pseudo-random bit sequence that it produces?

SOLUTION: Figure 11.55 shows an 8-bit LFSR using the four taps after the 1st, 6th, 7th, and 8th bits, as given in Table 11.7. It produces a sequence of $2^8 - 1 = 255$ bits before repeating.

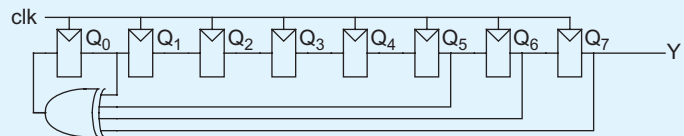


FIGURE 11.55 8-bit LFSR

$$\begin{array}{ll}
 \text{Binary} \rightarrow \text{Gray} & \text{Gray} \rightarrow \text{Binary} \\
 G_{N-1} = B_{N-1} & B_{N-1} = G_{N-1} \\
 G_i = B_{i+1} \oplus B_i & B_i = B_{i+1} \oplus G_i \quad N-1 > i \geq 0
 \end{array} \quad (11.30)$$

11.7.4 XOR/XNOR Circuit Forms

One of the chronic difficulties in CMOS circuit design is to construct a fast, compact, low-power XOR or XNOR gate. Figure 11.59 shows a number of common static single-rail 2-input XOR designs; XNOR designs are similar. Figure 11.59(a) and Figure 11.59(b) show gate-level implementations; the first is cute, but the second is slightly more efficient. Figure 11.59(c) shows a complementary CMOS gate. Figure 11.59(d) improves the gate by optimizing out two contacts and is a commonly used standard cell design. Figure 11.59(e) shows a transmission gate design. Figure 11.59(f) is the 6-transistor “invertible inverter” design. When A is 0, the transmission gate turns on and B is passed to the output. When A is 1, the A input powers a pair of transistors that invert B . It is compact, but nonrestoring. Some switch-level simulators such as IRSIM cannot handle this unconventional design. Figure 11.59(g) [Wang94] is a compact and fast 4-transistor pass-gate design, but does not swing rail to rail.

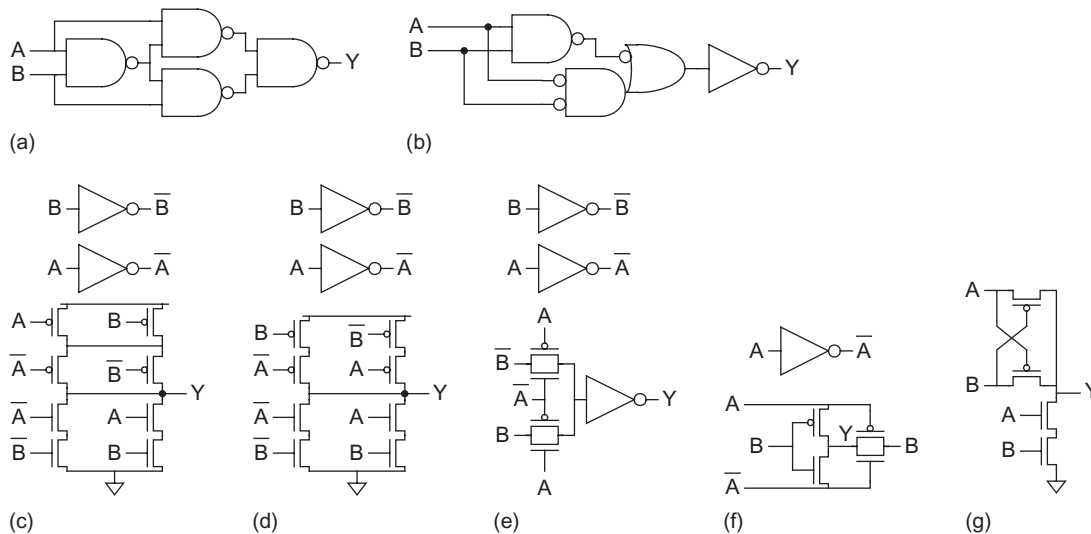


FIGURE 11.59 Static 2-input XOR designs

XOR gates with 3 or 4 inputs can be more compact, although not necessarily faster than a cascade of 2-input gates. Figure 11.60(a) is a 4-input static CMOS XOR [Griffin83] and Figure 11.60(b) is a 4-input CPL XOR/XNOR, while Figure 9.20(c) showed a 4-input CVSL XOR/XNOR. Observe that the true and complementary trees share most of the transistors. As mentioned in Chapter 9, CPL does not perform well at low voltage.

Dynamic XORs pose a problem because both true and complementary inputs are required, violating the monotonicity rule. The common solutions mentioned in Section 11.2.2.11 are to either push the XOR to the end of a chain of domino logic and build it

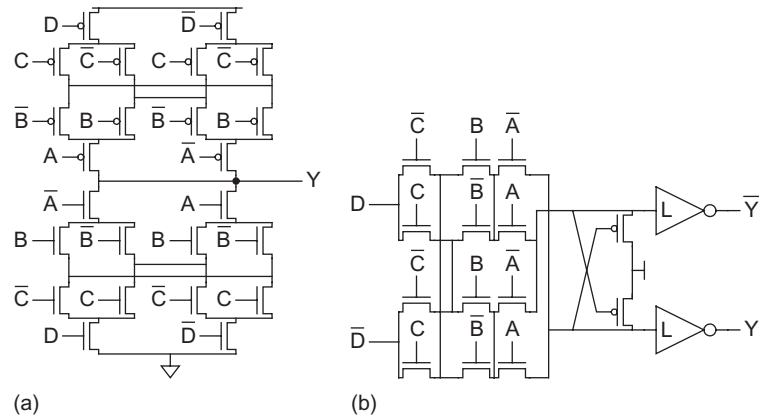


FIGURE 11.60 4-input XOR designs

with static CMOS or to construct a dual-rail domino structure. A dual-rail domino 2-input XOR was shown in Figure 9.30(c).

11.8 Shifters

Shifts can either be performed by a constant or variable amount. Constant shifts are trivial in hardware, requiring only wires. They are also an efficient way to perform multiplication or division by powers of two. A variable shifter takes an N -bit input, A , a shift amount, k , and control signals indicating the shift type and direction. It produces an N -bit output, Y . There are three common types of variable shifts, each of which can be to the left or right:

- **Rotate:** Rotate numbers in a circle such that empty spots are filled with bits shifted off the other end
 - Example: 1011 ROR 1 = 1101; 1011 ROL 1 = 0111
- **Logical shift:** Shift the number to the left or right and fills empty spots with zeros.
 - Example: 1011 LSR 1 = 0101; 1011 LSL 1 = 0110
- **Arithmetic shift:** Same as logical shifter, but on right shifts fills the most significant bits with copies of the sign bit (to properly sign, extend two's complement numbers when using right shift by k for division by 2^k).
 - Example: 1011 ASR 1 = 1101; 1011 ASL 1 = 0110

Conceptually, rotation involves an array of N N -input multiplexers to select each of the outputs from each of the possible input positions. This is called an *array shifter*. The array shifter requires a decoder to produce the 1-of- N -hot shift amount. In practice, multiplexers with more than 4–8 inputs have excessive parasitic capacitance, so they are faster to construct from $\log_v N$ levels of v -input multiplexers. This is called a *logarithmic shifter*. For example, in a radix-2 logarithmic shifter, the first level shifts by $N/2$, the second by $N/4$, and so forth until the final level shifts by 1. In a logarithmic shifter, no decoder is necessary. The CMOS transmission gate multiplexer of Figure 9.47 is especially well-suited to logarithmic shifters because the hefty wire capacitance is driven directly by an inverter rather than through a pair of series transistors. 4:1 or 8:1 transmission gate multiplexers reduce the number of levels by a factor of 2 or 3 at the expense of more wiring and