

Experiment 3

Nonlinear circuits: diodes and analog multipliers

So far the analog circuits we have considered have all been linear, so that the output has been given by a sum of terms, each term strictly proportional to only one source value (see the section in Experiment 1: *Linear circuits and superposition* on page 1-19, and, in particular, equation 1.8 on 1-20). It's now time to extend our design toolbox to include *nonlinear* elements and networks, ones for which the *principle of linear superposition* no longer holds, or, often, holds only partly or under more restrictive conditions.

One of the simplest nonlinear circuit components is the *semiconductor diode*, which we consider first. This two-terminal element behaves in a most asymmetric manner: its resistance is very low for currents of more than a few milliamps flowing in one direction through the device, but it has an enormously high resistance to current flow in the opposite direction. This element is very useful for constructing absolute value, peak detection, overvoltage protection, and more general nonlinear resistance circuits. The diode's current-voltage relationship is actually exponential, so it is also useful for building exponential and logarithmic response amplifiers. Its characteristics are strongly temperature-dependent, so a diode also makes an excellent, accurate temperature sensor. Of course, some types of diodes also can emit and detect light (LEDs, laser diodes, photodiodes and solar cells), so the variety of applications of the seemingly simple semiconductor diode is nearly endless.

The other nonlinear element we'll consider in this experiment is the much more sophisticated *analog multiplier* integrated circuit, whose output is proportional to the *product* of two input voltages (making its transfer function a so-called *bilinear form* of its two inputs). This flexible device may be used in circuits which not only multiply but also divide, raise to powers, and take roots. You may use it to build variable-gain amplifiers, modulators and demodulators, phase detectors, automatic gain control and signal compression circuits, voltage-controlled filters, and oscillators — as well as its obvious applications in general *analog computing* circuits.

Copyright © Frank Rice 2013, 2019
Pasadena, CA, USA
All rights reserved.

CONTENTS

TABLE OF CIRCUITS	3-IV
THE DIODE AS A RECTIFIER	3-1
<i>A simple semiconductor diode model</i>	3-1
<i>Basic half-wave rectifier circuit</i>	3-2
<i>Precision rectifier circuits</i>	3-3
<i>Peak detectors and AM demodulation</i>	3-6
THE ANALOG MULTIPLIER	3-9
<i>The ideal analog multiplier</i>	3-9
<i>A real analog multiplier IC</i>	3-11
PRELAB EXERCISES	3-13
LAB PROCEDURE	3-15
<i>Overview</i>	3-15
<i>Diode half-wave rectifiers</i>	3-15
<i>AM demodulator</i>	3-16
<i>Using an analog multiplier as a frequency doubler</i>	3-16
<i>Amplifier with voltage-controlled gain</i>	3-17
<i>Additional, self-directed investigations</i>	3-17
<i>Lab results write-up</i>	3-17
MORE CIRCUIT IDEAS	3-18
<i>Exponential and logarithmic amplifiers</i>	3-18
<i>A fast peak detector</i>	3-19
<i>True RMS measurement using analog multipliers</i>	3-21
ADDITIONAL INFORMATION ABOUT THE TEXT IDEAS AND CIRCUITS	3-23
<i>Zener diode regulator</i>	3-23
<i>LEDs</i>	3-25
<i>Using a PN junction for temperature sensing</i>	3-27
<i>Approximating a transcendental function using analog multipliers</i>	3-28
<i>Full-wave and bridge rectifiers</i>	3-31
THE PHYSICS OF THE PN JUNCTION DIODE	3-35
<i>Insulators, conductors, and semiconductors</i>	3-35
<i>Electrons and holes; impurities and doping</i>	3-37
<i>The equilibrium PN junction</i>	3-40
<i>The PN junction I-V characteristic curve</i>	3-41
<i>Zener and avalanche breakdown</i>	3-43

TABLE OF CIRCUITS

Rectifier, half-wave, simple	3-2
Rectifier, half-wave, precision	3-4
Rectifier, full-wave, precision	3-4
Rectifier, full-wave, transformer-driven	3-32, 3-33
Rectifier, bridge, transformer-driven	3-34
Rectifier, half-wave, capacitor filtered	3-6
Peak detector, precision	3-7
Peak detector, precision, fast response	3-20
Absolute value circuit	3-4
Analog multiplier functional circuit, MPY634	3-11
Analog divider circuit	3-9, 3-12
Analog square root circuit	3-9, 3-12
Triangle-sine converter	3-31
Zener diode voltage regulator	3-24
LED pilot light circuits	3-27
Temperature sensor, PN diode	3-28

THE DIODE AS A RECTIFIER

A simple semiconductor diode model

A semiconductor *diode* is a two-terminal element which acts as a “one-way valve” for electrical current (i.e., it is a *rectifier*). The most common type of diode is made from a silicon crystal divided into two layers with differing impurity atoms mixed into the silicon. The resulting structure creates a *PN junction* at the interface between the two layers which gives the diode its rectification property. A simplistic, mostly qualitative description of this rather complicated phenomenon of solid-state physics is provided in the section *The physics of the PN junction diode* starting on page 3-35.

The schematic symbol of a diode is shown at right along with a photo of a typical silicon signal diode. When the diode is *forward-biased* its resistance becomes very small, and current will flow through it in the direction shown (note that the schematic symbol includes an arrow (triangle) which points in the direction of the current flow). Forward-biasing is accomplished by applying a voltage so that the diode’s *anode* is at a more positive (+) voltage than its *cathode*. As you can see in the figure, *the cathode is denoted by a line in the schematic symbol and is usually marked by a line or stripe on the physical diode’s body*.

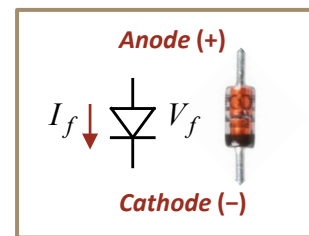


Figure 3-1: A typical silicon signal diode (type 1N914 or 1N4148) and its schematic symbol. This diode’s glass case is actually only about 3mm long.

When a diode is forward-biased and conducts a current of at least a few milliamps, the voltage drop across its two terminals (V_f in Figure 3-1) remains nearly constant even for an increase in the forward current (I_f) of an order of magnitude or more. When a diode is *reverse-biased* (anode more negative than the cathode), on the other hand, only a very small *leakage current* flows through it. This leakage current is quite insensitive to the reverse voltage applied to the diode, at least until some critical reverse voltage is reached which causes the diode to suddenly *break down* (and, usually, catastrophically fail). Therefore we can construct the following first approximation of its characteristics (good enough to use for many applications):

LOWEST-ORDER DIODE CHARACTERIZATION

A diode’s basic behavior is characterized by the following two parameters:

- V_f **forward voltage drop:** the nearly constant voltage drop across the diode as it conducts current (of some specified magnitude) in its forward-biased direction.
- I_R **reverse leakage current:** the small current which flows through the diode when the applied voltage is less than V_f or whenever the diode is reverse-biased.

Experiment 3: The diode as a rectifier

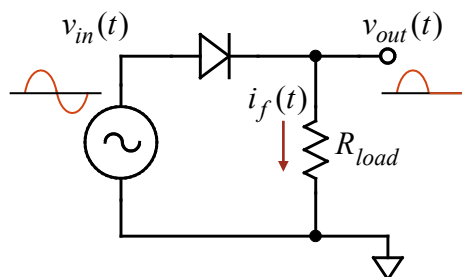
The forward voltage drop, V_f , of a PN junction diode is determined by the semiconductor material from which it is constructed; in the case of silicon (Si) conducting a few milliamps of current this voltage is approximately 0.5–0.6V, whereas for germanium (Ge) V_f is 0.3–0.4V, and for LEDs V_f is greater than 1.7V. The reverse leakage current, I_R , may vary by a couple of orders of magnitude depending on the type of diode and its temperature, but is generally a few tenths of a microamp or less (and may be much less). For a *photodiode* this reverse leakage current increases linearly with the intensity of any light shining on it. A *perfect diode* would be one with both $V_f \equiv 0$ and $I_R \equiv 0$.

A *forward-biased* diode is said to be **on** or **conducting**. A *reverse-biased* diode is said to be **off**.

Basic half-wave rectifier circuit

A very common diode rectifier circuit is shown in Figure 3-2; we will analyze it using our simple diode model. The input source to the circuit provides an AC signal $v_{in}(t)$; the circuit *rectifies* the input to produce an output $v_{out}(t)$ which is always of the same polarity across the load R_{load} (in this case, $v_{out}(t) \geq 0$).

Figure 3-2: Diode half-wave rectifier circuit. The AC driving source is converted to a rectified output (always of the same polarity); with the chosen orientation of the diode $v_{out}(t) \geq 0$. As shown by the waveform plots, only the positive half of an AC input cycle gets to the output.



The half-wave rectifier circuit is simple to analyze, particularly if we make the assumption that the diode is perfect (both V_f and $I_R \equiv 0$). First note that the diode's cathode is connected to ground (and the bottom terminal of the AC source) through R_{load} . This connection ensures that *each terminal of the diode always has a well-defined voltage* even when the diode is not conducting (turned off). It is important that your circuit designs follow this rule.

Whenever $v_{in}(t) > 0$, the diode's anode will be more positive than its cathode, and if $v_{in}(t) \geq V_f$ the diode will conduct current. In this case current will flow through R_{load} and return to the source. Assuming $V_f = 0$ implies that $v_{out} = v_{in}$. When $v_{in}(t) < 0$, on the other hand, the polarity across the diode reverses (because its cathode is still connected to ground through R_{load}). The diode turns off, and, assuming $I_R = 0$, then $v_{out} = 0$. So the output waveform will appear as in the figure: only the positive half of an input cycle reaches the output. Conversely, if the diode's orientation were reversed, then the output waveform would include only the negative half of an input cycle.

DIODE OPERATING STATES

We analyzed the operation of the half-wave rectifier circuit by *separately considering the two diode operating states*: conducting and turned off. We can then separately analyze the behavior of the circuit for each diode operating state using all of our standard, linear methods. As part of this analysis we can also determine the ranges of input and output signal values which will place the diode into one or the other state. This is the basic approach we will use to analyze nearly all of our diode circuits.

If a diode is not perfect (and no diode is), then the diode remains off until the voltage drop across it has the correct polarity and reaches its forward voltage drop V_f . Thus V_f will be lost across it when it does conduct (which is why V_f is also commonly called a *diode drop*). Therefore in the case of the half-wave rectifier in Figure 3-2, when the diode is on $v_{out} = v_{in} - V_f$. (see the figure at right). Similarly, a nonzero I_R will flow in the opposite direction through the diode when it is reversed biased (turned off); for the half-wave rectifier this implies that v_{out} can become very slightly negative: $v_{out} = -R_{load} I_R$.

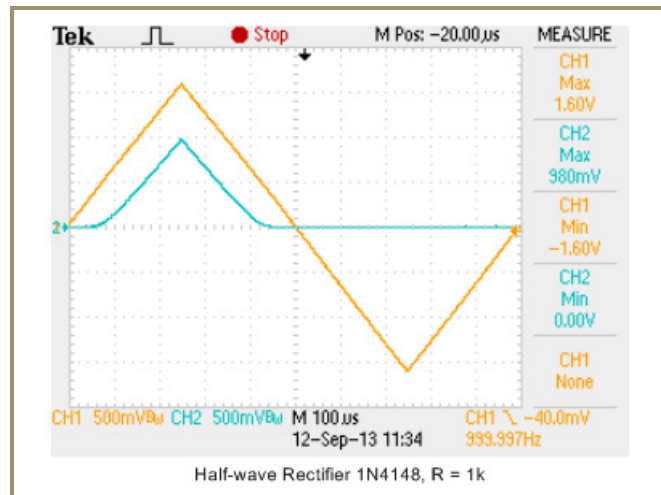


Figure 3-3: Response of the half-wave rectifier showing the effect of V_f , the diode's forward voltage drop. The diode is silicon, so $V_f \approx 0.6\text{V}$ (input is CH1, output CH2). Note that the diode actually turns on gradually, as indicated by the curve in the trace for v_{out} near 0.

Precision rectifier circuits

A diode's forward voltage drop of a few tenths of a volt (Figure 3-3) means that the basic half-wave rectifier circuit of the previous section can hardly be considered to provide precision rectification of an input signal, especially if the signal is small. As you might expect, adding the capabilities of an operational amplifier can remedy this situation: consider the right-hand circuit in Figure 3-4 on page 3-4.

This circuit is essentially a voltage follower, but a diode has been added in series with the op-amp output *inside the feedback loop*. When $v_{in}(t) > 0$, the op-amp will set its output one diode drop higher than v_{in} so that its *-Input* will equal v_{in} . Thus, for $v_{in}(t) > 0$, $v_{out} = v_{in}$ even though the diode may have a significant forward voltage drop. When $v_{in}(t) < 0$, the op-amp's output will go negative, and, because the diode's cathode is connected to ground through

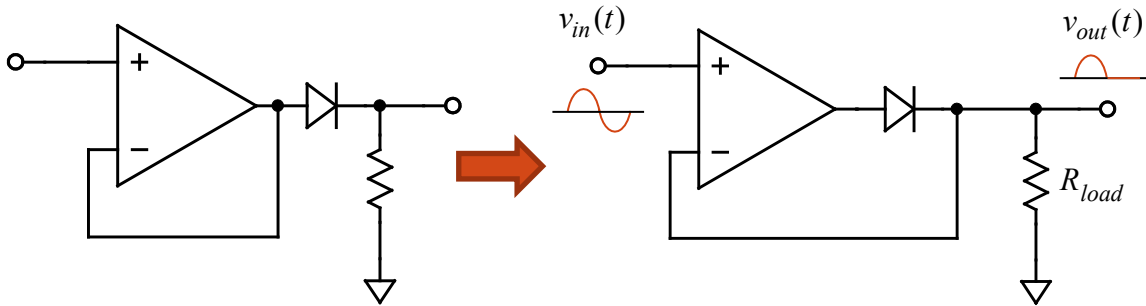


Figure 3-4: Simple, precision half-wave rectifier. Starting from a simple rectifier driven by a voltage follower, its feedback connection is moved to the other terminal of the diode. Placing the diode inside the feedback loop ensures that the op-amp output will compensate for the diode's forward voltage drop. When $v_{in} < 0$, the diode turns off as the op-amp output goes negative; R_{load} then ensures that $v_{out} = 0$ in this case.

R_{load} , the diode will turn off, disconnecting the op-amp's output from v_{out} . Thus, for $v_{in}(t) \leq 0$, $v_{out} = 0$. The load resistor value should usually be between 1k and 100k; if low output impedance is required even when the output is 0, buffer the output with a voltage follower. Note that when $v_{in} < 0$ the op-amp +Input voltage will be less than its -Input voltage (which will then be 0), so the op-amp's output will be at negative saturation, about -11V for the TL082 with $\pm 12V$ power supplies. When $v_{in}(t)$ rises back up through 0 the op-amp's output will quickly slew from negative saturation back up through 0 so that it will forward bias the diode and recapture the condition $v_+ = v_-$. An example of the output from this circuit is shown in the left-hand image in Figure 3-6 on page 3-5.

By adding another op-amp we can construct a precision *full-wave rectifier*, or *absolute value circuit*, Figure 3-5. The second amplifier (using op-amp U2) combines the original input signal with the half-wave rectifier output to form the absolute value of the input waveform. The two resistors used to form the feedback network for U2 should be well-matched in value. Note that U2 and the two resistors R form the combination inverting+noninverting amplifier

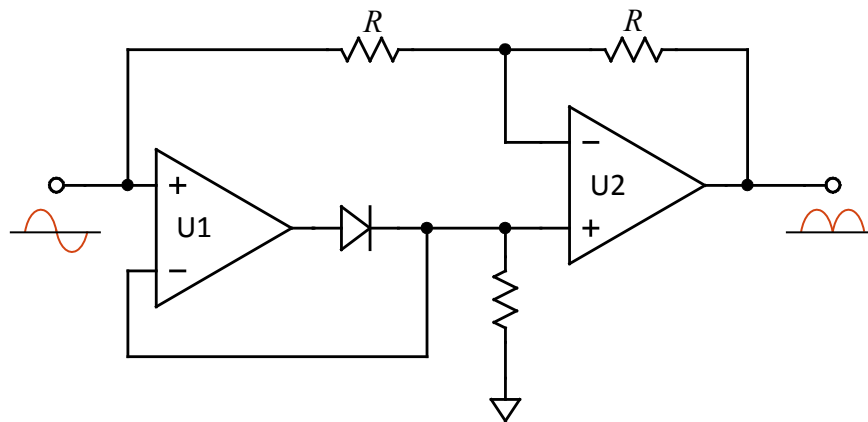


Figure 3-5: Precision full-wave rectifier (absolute value circuit). The output of a half-wave rectifier is combined with the original input signal by the amplifier U2. Its output is the absolute value of the input. The two resistors R should be well-matched in value (1% or better).

circuit we saw back in Experiment 1 (Figure 1-22 on page 1-24). Its v_{in-} input is just the input signal to be rectified, whereas its v_{in+} input comes from the output of a half-wave rectifier circuit formed using U1. Analysis of this circuit is left to the exercises. An example of the output from this circuit is shown in the right-hand image in Figure 3-6.

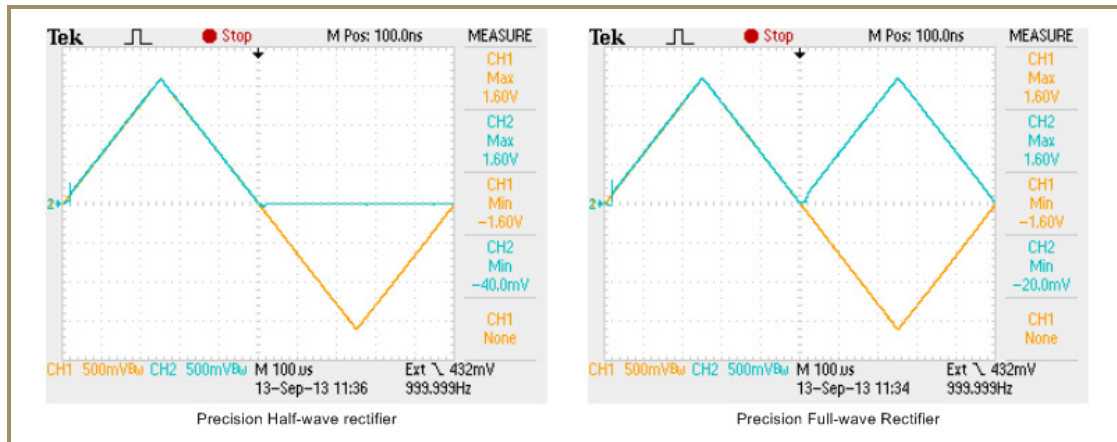


Figure 3-6: Precision rectifier outputs. Left: half-wave rectifier circuit of Figure 3-4. Right: full-wave rectifier of Figure 3-5. Input signal waveform and display scales are the same as in Figure 3-3 on page 3-3; note that the precision rectifier circuits eliminate the 0.6V forward voltage drop of the silicon diode used in the circuits. The small “glitches” in the output waveforms as the input goes through 0 are effects of the TL082 op-amp’s finite slew rate (explored further in Experiment 4).

Before proceeding further, some important diode limitations should be noted. For the 1N4148 silicon small-signal diodes (the type you will mostly use) the limits are:

$$V_R = 75\text{ V} \quad I_F = 300\text{ mA} \quad P_D = 0.3\text{ W}$$

DIODE LIMITATIONS

Exceeding these limits could cause catastrophic diode and circuit failure:

V_R reverse breakdown voltage: the maximum reverse-bias voltage which may be safely applied without the diode exhibiting avalanche or Zener breakdown.

I_F maximum forward current: the maximum current the diode can tolerate when forward-biased.

P_D maximum power dissipation: the maximum power dissipation the diode can tolerate without overheating and failing.

Warning

Because the forward voltage drop, V_f , of a diode is nearly independent of the current flowing through it, *the magnitude of the forward current must be limited by the external circuit*, or a forward-biased diode will quickly fail.

Peak detectors and AM demodulation

If you put a capacitor in parallel with the output of a simple, half-wave rectifier as shown in Figure 3-7 below, then whenever the diode conducts the AC voltage source will quickly charge up the capacitor to match its voltage (minus the diode's forward voltage drop). When the AC voltage source passes its peak positive output and starts to decrease, its voltage drops below the capacitor's voltage, and the diode turns off. The capacitor then begins to relatively slowly discharge through the load resistor until the AC source voltage output again becomes high enough to turn the diode back on and recharge the capacitor. As a result, the output voltage stays near the peak positive input voltage value; the smaller the capacitor or the smaller the load resistor, then the more the capacitor will discharge during a cycle. This simple circuit is useful as part of a *power supply* to turn the 60Hz power line AC voltage (usually coupled through a transformer) into a nearly constant DC voltage to power an electronic device. The amplitude of the capacitor charge-discharge output oscillation is called the power supply *ripple voltage*; small ripple voltage is desirable, so power supply filter capacitors tend to be large, and they are usually electrolytic (refer again to the photo in Experiment 2, Figure 2-1 on page 2-6).

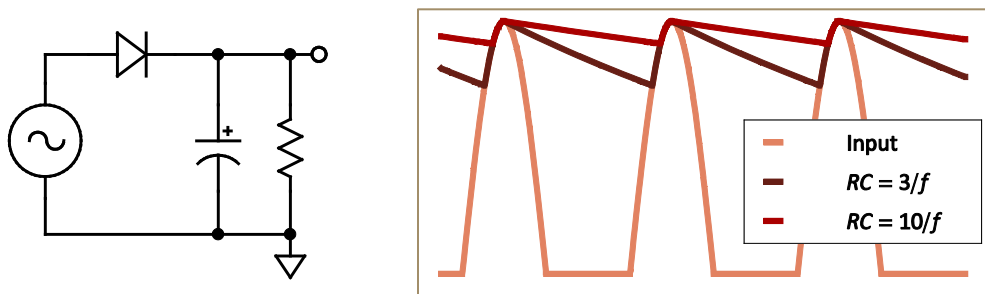


Figure 3-7: Using a filter capacitor to smooth the output of a simple, half-wave rectifier. The capacitor C is charged whenever the input voltage exceeds the output voltage so that the diode is forward-biased; when the input voltage drops and the diode becomes reverse-biased, the capacitor discharges through the load (the resistor R). The graph at right shows the resulting output voltage variation for various RC time constant values; f is the frequency of the input source. The effect of the diode's forward voltage drop is not included in the graph.

If we use this idea with a precision half-wave rectifier circuit, we get the *peak detector* shown in Figure 3-8, which works just the same as the simple circuit discussed above, except now an op-amp is used to correct for the diode forward voltage drop and a voltage follower is added so that the subsequent load will not discharge the capacitor. This circuit can hold the maximum positive input voltage encountered for quite a long time; if the resistor R is removed and a high-quality capacitor is used, then the capacitor's discharge will only be because of the follower op-amp's input bias current (typically 50 pA for the TL082 at room temperature) and the diode's reverse leakage current (25 nA max for the 1N4148 at room temperature), the sum of which could be kept to only a few nanoamps by careful component selection. In this case, a 10 μ F capacitor would provide an output voltage which decayed at less than 1 mV/sec. If you need to hold an accurate peak voltage value for longer than this,

then your best bet would be to record a *digitized* version of the peak detector's output voltage value. Note that the half-wave rectifier op-amp's output will go to negative saturation when its input drops below the capacitor's voltage, because its *-Input* voltage will then be greater than its *+Input* voltage. This transition could be used to signal an *analog-to-digital converter* circuit to digitize and save the peak detector's output voltage.

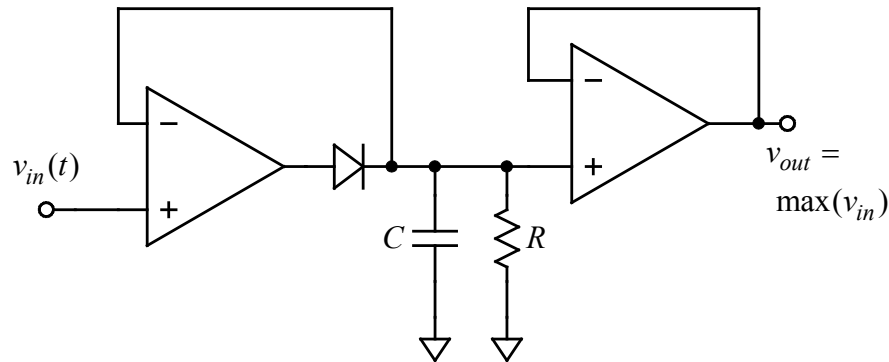


Figure 3-8: A peak detector circuit, which outputs the most positive value of the input signal seen so far; reversing the diode would output the most negative value of the input. The output voltage will exponentially decay toward 0 with time constant RC ; because of the voltage follower on the output, RC may be made quite long by choosing a large value for R , or the resistor could even be replaced by a switch to reset the output to 0.

AM demodulation

The earliest method used to transmit and receive voice (audio) communications at radio frequencies (~ 1 MHz and up) was to vary the amplitude of a radio-frequency (RF) *carrier* waveform in synchronization with the audio signal to be transmitted, a technique known as *amplitude-modulation* (AM). The idea is illustrated in Figure 3-9. To recover the signal a receiver is tuned to the RF carrier frequency. The information in the RF signal is then

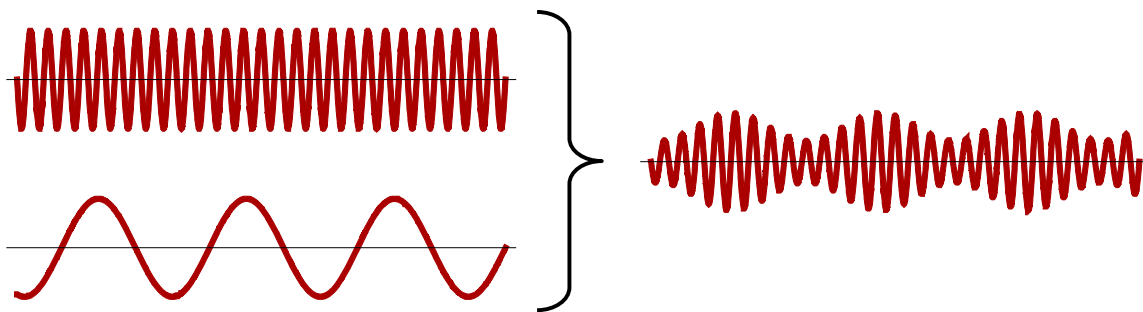


Figure 3-9: AM modulation. A high-frequency carrier wave has its amplitude varied in synchronization with a much lower frequency signal. The varying amplitude of the carrier then contains the information in the transmitted signal.

Experiment 3: The diode as a rectifier

recovered by a circuit (the *detector*) which tracks the RF carrier's amplitude in a process called *AM demodulation*.

The basic process used to demodulate the RF waveform is illustrated in Figure 3-10. For example, with a clever choice for the *RC* time constant used in our precision peak detector shown in Figure 3-8 we can use it to perform this function. If the *RC* time constant is long compared to the RF carrier wave's period, then the peak detector will output the carrier's amplitude. If at the same time the *RC* time constant is short compared to the AM signal's period, then the peak detector's output can decay fast enough to follow the changes in the carrier's amplitude. Thus the output will follow the AM modulated waveform's amplitude envelope as shown in the left-hand diagram in Figure 3-10. AC coupling this output into a subsequent circuit would remove the envelope's DC offset, recovering the original signal used to AM modulate the carrier wave.

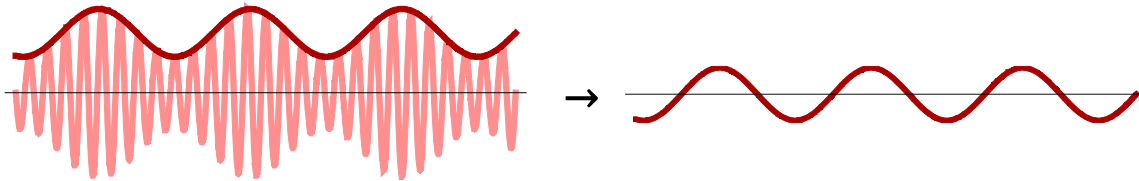


Figure 3-10: AM demodulation. The varying amplitude of the AM modulated carrier wave is followed by a detector which outputs the envelope of the wave. This output is then band-pass filtered to remove its DC offset and any residual RF carrier. The result is a recovered audio signal replicating the source information.

Actually, using a precision peak detector is, of course, unnecessarily complicated for simple AM detection. For example, a purely passive (no power supply) *crystal radio* to receive AM modulated radio transmissions was demonstrated as early as 1906 by the American engineer Greenleaf Whittier Pickard. His receiver consisted of just a simple, half-wave rectifier fed by a passive, *LC resonant circuit* tuned to the RF carrier wave's frequency.¹ His rectifying element was a metal-semiconductor junction whose modern version is called a *Schottky diode* (after German physicist Walter H. Schottky).

¹ We'll investigate resonant circuits in Experiment 5.

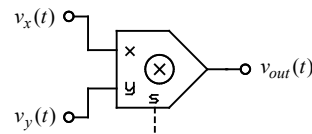
THE ANALOG MULTIPLIER

The ideal analog multiplier

Unlike the simple semiconductor diode, the modern analog multiplier is a sophisticated, complicated integrated circuit incorporating temperature-compensated voltage references, matched transistors, and laser-trimmed resistors. Much like the modern operational amplifier, this internal circuitry makes these devices particularly simple to utilize (although actual multipliers tend to behave in a less ideal manner than their op-amp cousins). In this section we consider the properties of an ideal analog multiplier and discuss applications of such a device; the next section will look at the use and limitations of an actual multiplier IC.

An analog multiplier, of course, requires at least two user-controlled input signals, which are typically called the x and y inputs. The input and output analog values may be voltages or currents, depending on the particular device, but modern medium-speed, precision multipliers usually input and output *voltages*, so that is what we will assume here. The product of two voltages has, of course, units of Volts², but the output will be in Volts. Thus the product xy is divided by an additional voltage parameter, the *scale factor*, s , which converts the product from Volts² to Volts (equation 3.1).

3.1 Multiplier:
$$v_{out} = v_x v_y / s$$



The analog multiplier usually has a precision, built-in scale factor which is often $s = 10.0\text{ V}$, but it may include an external terminal so that you can adjust the value of s (as shown in the schematic above). However the scale factor may be set, it is usually limited to positive values ($s > 0$). Modern analog multipliers don't limit the signs of v_x and v_y (they support “four-quadrant multiplication”), but they usually limit their magnitudes to less than s .

Placing a multiplier in an op-amp's feedback loop provides for the calculation of the inverse operation, division, as shown in Figure 3-11. Assuming the negative feedback around the op-

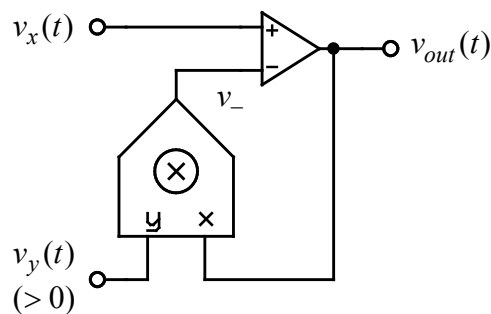


Figure 3-11: Basic divider circuit using an op-amp. Note that for the negative feedback to be effective, v_- must have the same sign as v_{out} , and this condition requires that $v_y > 0$. Improving this circuit to allow “four-quadrant division” so that $v_y < 0$ is supported is a nontrivial exercise and could make a good final project.

amp is effective, then $v_- = v_+ = v_x$. But the multiplier output is $v_- = v_{out}v_y/s$, so the divider circuit output is:

3.2 Divider: $v_{out} = s v_x / v_y \quad (v_y > 0)$

Since a typical multiplier will require that its inputs not exceed s , the gain of the circuit in Figure 3-11 with respect to the op-amp's *+Input* (v_x) will be greater than 1, and thus the circuit's bandwidth will be less than the op-amp's gain-bandwidth product f_{BW} ; in fact, the bandwidth will be $f_{BW}(v_y/s)$.

The condition $v_y > 0$ in (3.2) comes from the requirement that v_- must have the same sign as v_{out} (see Figure 3-11), or the feedback will effectively change sign and become *positive*, which will cause the op-amp's output to proceed to one of its power supply limits (i.e. the output will *saturate*) until v_y becomes positive.

You may calculate the square of a signal by connecting a multiplier's inputs together; the inverse operation, a square root circuit, again may be built by putting the multiplier in an op-amp feedback loop, but now circuit *latch-up* is a real possibility: the op-amp output will saturate (in this case at its negative limit) and *stay there* no matter what the input does. Using a precision rectifier configuration (as in Figure 3-4) is a popular way to avoid this latch-up problem (Figure 3-12).

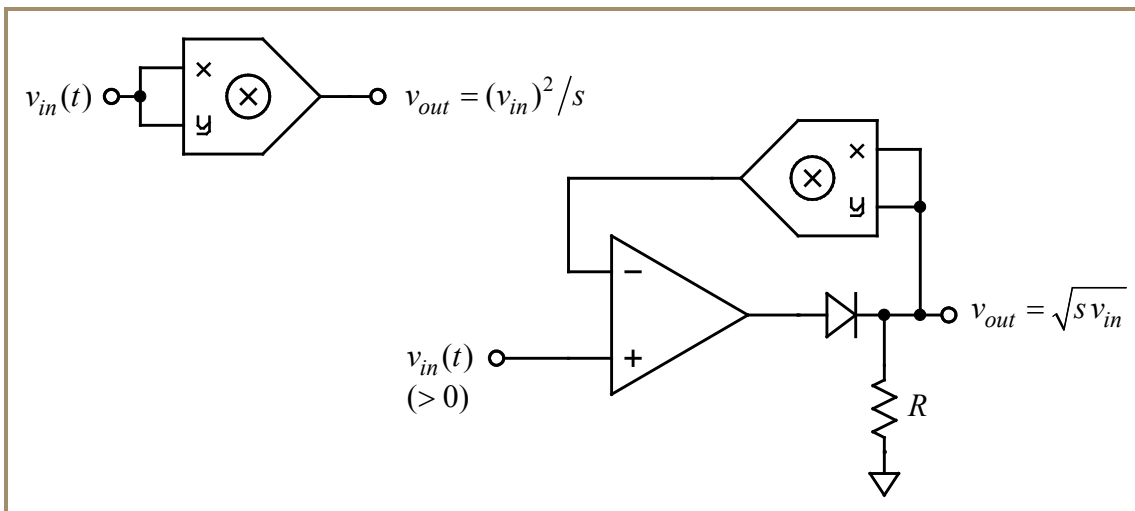


Figure 3-12: Square and square root circuits. The square root circuit is prone to *latch-up*: if the inputs to the multiplier were to go < 0 , even momentarily, its positive voltage output would drive the op-amp into saturation at its negative voltage limit, permanently maintaining this undesirable state. The diode+resistor combination on the op-amp output ensures that the multiplier inputs never go negative, avoiding latch-up. A value of about $1\text{k}\Omega - 10\text{k}\Omega$ is appropriate for R . Of course, s is the multiplier's *scale factor*.

A real analog multiplier IC

Because circuit designers often need to include op-amps and need to invert signals in their circuits, actual IC multiplier devices usually include extra circuitry to make the designer's job easier. In this section, we consider the [Texas Instruments MPY634](#) device, which is included on the *ASLK PRO* breadboard; a similar, general purpose, relatively inexpensive multiplier is the [Analog Devices AD633](#), which may be more readily available.² In this section we discuss the use of the quite complicated MPY634, a functional block diagram of which is shown below.

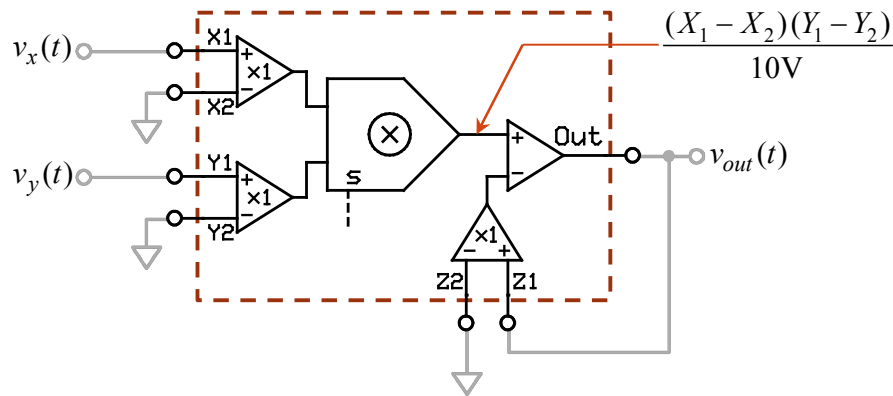


Figure 3-13: Functional block diagram of the Texas Instruments MPY634 analog multiplier. Besides the basic multiplier, the IC includes a general-purpose op-amp on the multiplier's output and has differential inputs for all parameters, including the Z inputs to the op-amp's feedback (-Input) terminal. Shown in light gray are the connections needed to emulate the basic, generic multiplier discussed in the previous section. Note that the $\times 1$ amplifiers on the multiplier inputs are not op-amps – they simply apply the difference of their two inputs to their outputs ($X_1 - X_2$, etc.).

The MPY634 has differential inputs for all input parameters and includes a general-purpose op-amp on the multiplier output so that it can be easily configured for different scale factors or for inverse operations. In Figure 3-13 the output op-amp has been configured as a voltage follower (output fed back through its Z1 input), and the negative differential input for each parameter has been grounded; the resulting circuit acts as the basic multiplier described by equation 3.1 in the previous section. Note that the default scale factor is 10V; this is the value if the IC's scale factor terminal (S in the figure) is left disconnected. The MPY634 basic multiplier accuracy in this configuration is approximately 2%. The MPY634 output op-amp's $f_{BW} \geq 6$ MHz, and it has a 20 V/ μ s slew rate, so it has similar performance to the TL082 op-amps on the breadboard.

Divider and square root configurations are shown in Figure 3-14 (on page 3-12). Because the multiplier output is connected to the internal op-amp's +Input, you have to *invert the multiplier's output* to provide negative feedback around the op-amp (as in Figure 3-11 and

² “Inexpensive” is a relative term. The TL082 IC you've used costs less than \$1 and includes two complete op-amp devices, whereas the AD633 IC cost over \$13 for a single multiplier, and the MPY634s on the *ASLK PRO* board cost nearly \$24 each (prices as of fall 2019).

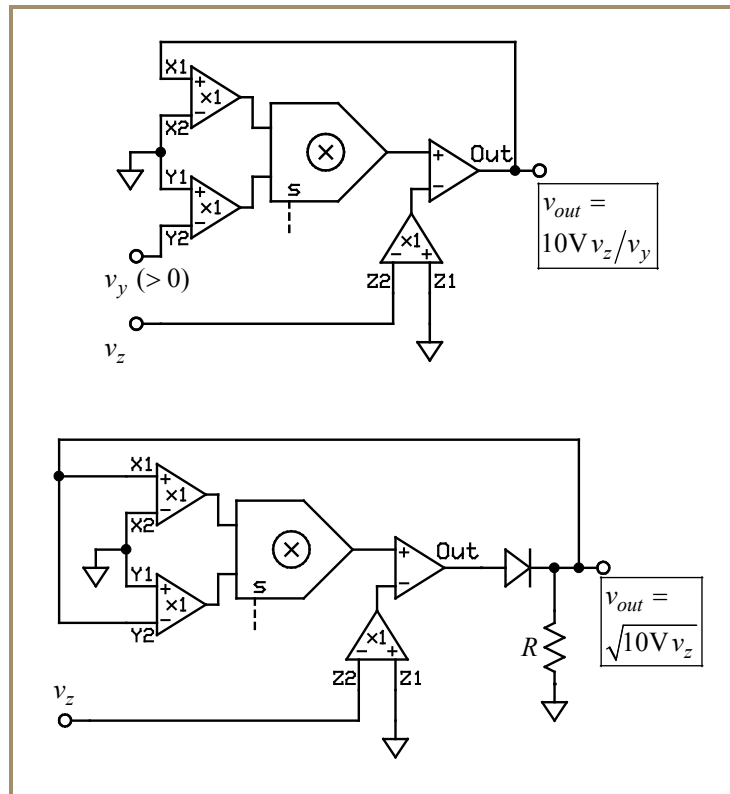
Experiment 3: The analog multiplier

Figure 3-12). This is accomplished by using the (-) differential terminal for one of the arguments to the multiplier, so its output is the negative of the product of its two inputs. Similarly, the other op-amp input (to its *-Input* terminal) is also inverted by using the input's (-) differential terminal. The [MPY634 data sheet](#) has more examples demonstrating this device's flexibility.

Figure 3-14: MPY634 divider (top) and square root (bottom) circuits.

In each case the multiplier sub-circuit must be in the *negative feedback loop* of the op-amp (see Figure 3-11 and Figure 3-12); to accomplish this using the MPY634, the multiplier output must be the *negative* of the product of its inputs, which is accomplished by *inverting only one* of its inputs. The v_z input to the op-amp is also inverted, so that it effectively becomes a *noninverting* op-amp input ($-1 \times -1 = +1$).

A diode and resistor are still needed to prevent latch-up of the square root circuit output, as discussed in the previous section. The resistor value R should be about $1k\Omega - 10k\Omega$, as before (Figure 3-12).



You must always provide negative feedback for the output op-amp of the MPY634.

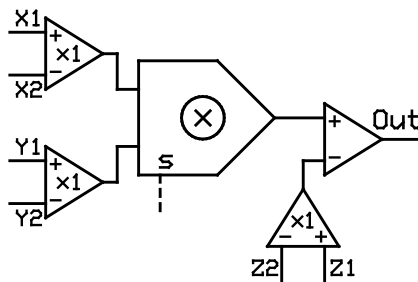
Any unused differential input terminals must be connected to ground.

PRELAB EXERCISES

1. Consider the simple half-wave rectifier circuit (left-hand schematic in Figure 3-2 on page 3-2). If the diode is perfect (both its forward voltage and reverse leakage current are 0), then what is the circuit's *output resistance* during a positive half-cycle (when $v_{in}(t) > 0$)? What about during a negative half-cycle ($v_{in}(t) < 0$)? Include R_{load} as part of the circuit when calculating the output resistance.

Hint: review the section on *output resistance* starting [on page 1-42 of Experiment 1](#). Consider two distinct cases: one when the diode is forward-biased (and, since $V_f \equiv 0$, its resistance = 0), and the other when the diode is reverse-biased (so, since $I_R \equiv 0$, its resistance is infinite).

2. How does the full-wave rectifier circuit in Figure 3-5 (on page 3-4) work? Analyze the two circuit operating states separately: (1) the diode is forward-biased; (2) the diode is reverse biased. Assume the op-amps are ideal; for each of these states determine:
 - a. the required input voltage (v_{in}) range for the circuit to be in that state
 - b. the voltage at U2's +Input
 - c. the circuit's transfer function v_{out}/v_{in}
3. Consider the square root circuit in Figure 3-12 on page 3-10. If the input $v_{in}(t) < 0$, then what is v_{out} ? What is the output voltage of the op-amp? What happens to the op-amp output and v_{out} as the input rises through 0? Does the forward-bias diode voltage drop affect the accuracy of v_{out} (assuming the op-amp is ideal)? Why or why not?
4. By configuring the MPY634's output op-amp to provide gain, you can effectively adjust the device's multiplication scale factor to something less than its 10V default. Complete the schematic diagram below by adding a feedback network to the op-amp and properly connecting the input terminals so that you have a multiplier with a scale factor $s = 5V$ (assign appropriate values to any resistors you include).



Another problem on the next page...

Experiment 3: Prelab exercises

5. Design a circuit using the MPY634 to create an amplifier with a *voltage-controlled gain*. The gain control input voltage range should be -7V to $+7\text{V}$ (minimum), and the circuit gain should equal the control voltage in volts, e.g. $-2\text{V} \rightarrow \text{gain} = -2$, etc. The allowable signal input and output voltage ranges (without *clipping* or distortion of the output) should be -7V to $+7\text{V}$ as well (as long as the circuit gain is not set too high). Use additional op-amp amplifier stages as part of the circuit design if you need them (but you may not need them). Assume that the input signal is ground-referenced, so you don't need a fully differential input for the signal. The gain accuracy should be no worse than $\pm 10\%$.

Provide a full schematic of the circuit with all component values included. Use only these standard resistor values: any of the following values times whatever power of 10 you need:

1.0, 1.1, 1.2, 1.3, 1.5, 1.6, 1.8, 2.0, 2.2, 2.4, 2.7, 3.0, 3.3, 3.6, 3.9, 4.3, 4.7, 5.1, 5.6, 6.2, 6.8, 7.5, 8.2, 9.1

For example, to get a resistor ratio of 5:1, you could use various combinations: 1.0×10^4 and 2.0×10^3 ohms, or 1.2×10^4 and 2.4×10^3 , or 1.5×10^4 and 3.0×10^3 , or 7.5×10^4 and 1.5×10^4 .

You will build and test this circuit during your lab session, so think carefully about it!

LAB PROCEDURE

Overview

The semiconductor diode and the analog multiplier are two very different examples of nonlinear components. Spend sufficient time during lab building circuits with each of them so that you become comfortable with using such elements in your designs. If you have time, combine circuits using them with the amplifier and filter designs you have already practiced with so that you start to think about and develop more complicated applications, such as those shown in later sections of this text.

Diode half-wave rectifiers

Build a simple half-wave rectifier circuit using a 10k load resistor and voltage followers to buffer the rectifier's input and output (Figure 3-15). Use bread-board op-amps which have preinstalled resistors and capacitors on their +Inputs to assemble the circuit, **but don't use a capacitor in the circuit yet**. Input a 1kHz triangle-wave signal with a couple of volts peak-peak amplitude and confirm that the diode's turn-on voltage is approximately 0.6V, as shown in Figure 3-3 on page 3-3. Next replace the diode with an LED. How to determine the LED polarity (anode v. cathode) is also shown in Figure 3-15. Note that you must increase the signal generator output to at least 4 or 5 V peak-peak to forward bias the LED. What is its approximate turn-on voltage?

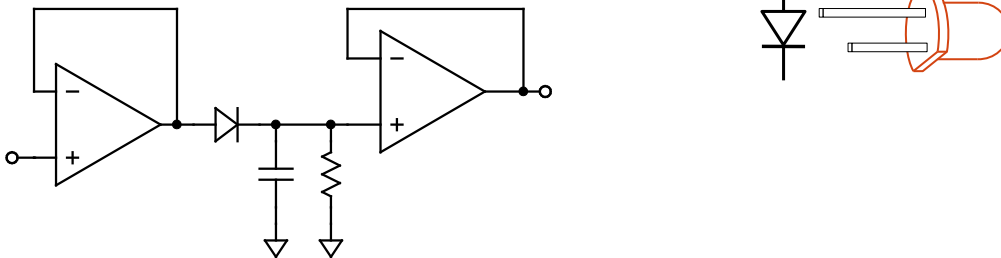


Figure 3-15 (left): Simple half-wave rectifier with voltage followers on its input and output. Use a resistor and capacitor already available on the +Input of the second op-amp for R and C . (right): How to determine the polarity of a LED.

Now move the feedback connection on the first op-amp to convert the circuit into a precision half-wave rectifier (Figure 3-4 on page 3-4). Note that the diode turn-on voltage drop is no longer evident in the circuit's output (compare your results to Figure 3-6 on page 3-5). Try it with both the LED and the silicon diode.

Now connect a capacitor to filter the rectified output as shown in Figure 3-15. First try the $0.1\mu\text{F}$ and then the $1\mu\text{F}$ capacitor available on the op-amp's +Input. Compare your results to the filtered half-wave rectifier output plotted in Figure 3-7 on page 3-6.

AM demodulator

Using the $0.1\mu\text{F}$ capacitor to filter your precision half-wave rectifier output, next configure the signal generator to produce an *amplitude modulated* (AM) signal. Use a sine-wave carrier frequency of 40kHz, a *modulation frequency* of 200Hz, and an AM *modulation depth* of 30%. Get your TA to help you set up the signal generator; trigger the oscilloscope using the rectifier circuit's output. Use about a 5V pk-pk carrier sine amplitude and adjust the scope's trigger level to steady its display. Your scope display should appear something like the screen capture image in Figure 3-16 below. Does the rectifier output detect the modulation (output mostly the modulated waveform's 200Hz *envelope*)? Try a couple of different load resistor and filter capacitor combinations. Vary the modulation frequency and modulation depth and note the effects on the output.

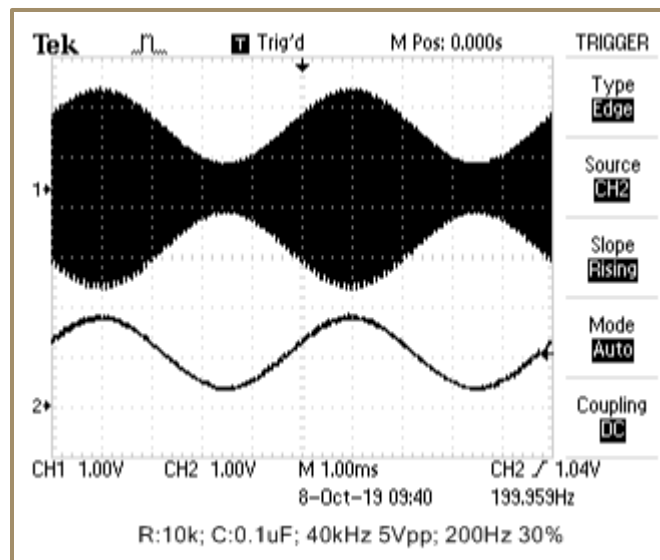


Figure 3-16: Oscilloscope screen capture of the AM demodulator input and output. The displayed menu shows the oscilloscope trigger setup. The component values and AM signal input parameters are shown in the image caption. A red LED was used as the rectifier.

Using an analog multiplier as a frequency doubler

Deselect the AM modulation on the signal generator and disassemble your rectifier circuit. Construct a basic squarer circuit using one of the analog trainer's three MPY634 multipliers; refer to Figure 3-13 on page 3-11 for the basic multiplier connections, and then connect the multiplier's X1 and Y1 inputs together and to the signal generator output. Using a 5V pk-pk sine-wave input signal at 1kHz, confirm that the squarer's output is given by $v_{out}(t) = v_{in}(t)^2 / 10\text{V}$ (what should be the output waveform's maximum and minimum voltages given the input you've provided?).

Once the circuit is working and behaves as expected, note that the squarer circuit's response to a sinusoid input is, of course, another sinusoid at twice the input frequency along with a

constant (DC) offset (since, of course, the output of the squarer is nonnegative). You may remove the DC component by *AC coupling* the output, as you learned in Experiment 2.

Assemble a gain 11 op-amp amplifier with an AC-coupled *RC* filter input (as in Experiment 2) and attach it to the multiplier's output. You have now constructed a *frequency doubler* — a sinusoid input produces a sinusoid output at twice the input frequency. Using the 1 kHz input frequency, determine what input signal amplitude is required to produce an output of the same amplitude. Increase the input frequency until you find the -3 dB upper bandwidth limit of your frequency doubler.

Amplifier with voltage-controlled gain

Build and test the circuit you designed for prelab exercise 5. The gain control voltage may be generated by a computer *DAQ* analog output port which you can control using the National Instruments *Measurement & Automation Explorer* application; the lab instructor or your TA will show you how to accomplish this.

Additional, self-directed investigations

Maybe try building a cubing circuit: $v_{out} \propto v_{in}^3$; how many multipliers would this take? Consider the full-wave rectifier circuit or one or more circuits from the **MORE CIRCUIT IDEAS** section (such as the true RMS circuit). Look back at earlier experiments to see if there are any other circuits you would like to investigate.

Lab results write-up

As always, include a sketch of the schematic with component values for each circuit you investigate, along with appropriate oscilloscope screen shots and, if appropriate, Bode plots. Make sure you've answered each of the questions posed in the above sections.

MORE CIRCUIT IDEAS

Exponential and logarithmic amplifiers

The exponential voltage-current relationship of the PN junction diode (equation 3.10 on page 3-42) may be exploited to build amplifiers with approximately exponential or logarithmic gains. If the diode current is much greater than I_R in equation 3.10, then the I-V relationship is approximately:

$$I = I_R e^{q_e V / \eta k_B T}$$

where, for a silicon diode, $I_R \sim 1\text{nA}$ and $q_e / \eta k_B T \approx 20 \text{ volt}^{-1}$ at room temperature. Thus, V will change by approximately 0.12V for a factor of 10 change in the current I , and $V \approx 0.6 \text{ V}$ for $I = 1\text{mA}$. Thus we can take the logarithm or the antilog (exponential) of an input voltage using the circuits in Figure 3-17 below. With resistor value R , the transfer functions are approximately:

$$|V_{out}| \approx 0.6\text{V} + 0.12\text{V} \times \log_{10} \left| \frac{V_{in}/\text{Volt}}{R/\text{k}\Omega} \right| \quad (\text{log amplifier})$$

3.3

$$|V_{out}| \approx (10^{-5} \text{ mA} \times R) 10^{|V_{in}/0.12\text{V}|} \quad (\text{exponential amplifier})$$

The sign of the output voltage in each case is the opposite of the sign of the input, since the amplifiers are inverting. The gains and offsets of these circuits may be inconvenient, so more circuitry is usually added to scale and offset the output signals; the transfer function, however, is quite temperature-dependent.

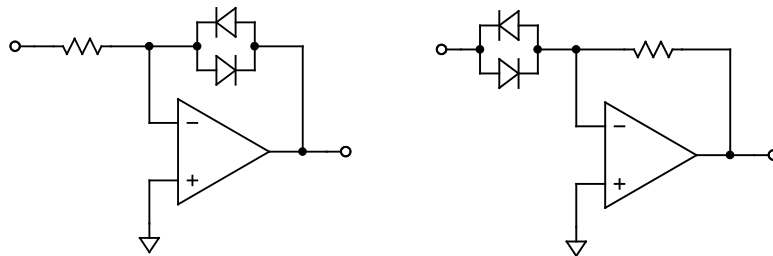


Figure 3-17: Inverting log (left) and exponential (right) circuits using diodes. These simple circuits are not very accurate and are very temperature-sensitive, but will work for noncritical applications. The paralleled back-to-back diodes will treat positive and negative input voltages the same, using whichever diode is forward-biased to dominate the circuit transfer function.

A fast peak detector

One potential problem with the peak detector shown Figure 3-8 on page 3-7 is that it can output a voltage significantly lower than the peak amplitude of a high frequency input signal. This is a big problem for some applications, so now consider ways to improve the peak detector's speed. There are two major design issues which increase the time it takes to change the capacitor's voltage: (1) the amount of current the op-amp output can supply, and (2) the op-amp *slew rate*. The first issue is easy to deal with, so we consider it first. When the input exceeds the voltage currently stored by the capacitor, the diode becomes forward-biased, and the op-amp output can charge the capacitor toward the new maximum voltage. The rate that the capacitor's voltage will change is given by $i_{max} = C dv/dt$, where i_{max} is the maximum output current the op-amp is capable of supplying (usually specified in the op-amp's data sheet). The TL082, for example, can output up to about 40mA into a discharged capacitor, but its available output current decreases as its output voltage rises. If your application must track a signal whose peak amplitude changes rapidly, use as small a capacitor value C as you can (consistent with your required peak hold time) or get an op-amp with a larger current output capacity.

The second problem is more difficult to handle. As we'll explore more thoroughly in Experiment 4, the op-amp has a maximum rate that it can change its output voltage. This output voltage *slew rate* is specified in the op-amp data sheet; the TL082, for example has a typical slew rate of 13 V/ μ s, but it could be as low as 8 V/ μ s. When used in an amplifier circuit configuration, an op-amp's slew rate limitation is usually compatible with its unity gain bandwidth and rarely causes a problem. Unfortunately, the design of the circuit in Figure 3-8 exacerbates the speed limitation imposed by the op-amp slew rate.

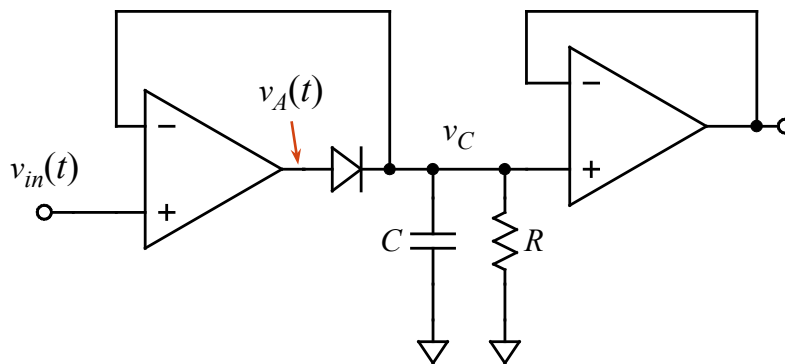


Figure 3-18: The peak detector of Figure 3-8, revisited.

Consider again that original peak detector circuit, repeated in Figure 3-18 above. Whenever the input voltage $v_{in}(t)$ is less than v_C , the voltage stored by the capacitor, the op-amp output voltage $v_A(t)$ goes all the way down to its negative limit (saturation), because the op-amp's *-Input* voltage will be greater than its *+Input* voltage. If and when the input returns to above

Experiment 3: More circuit ideas

v_C , the op-amp output must slew up from its negative limit (about -11V for the TL082) to a diode-drop above v_C in order to forward bias the diode and start to increase the capacitor's charge. The time it takes the op-amp to change its output voltage is, of course, limited by its slew rate. Since the output starts at the op-amp's negative saturation voltage, its output must change by several volts, which will take about a microsecond or more for the TL082. Clearly, this time delay will limit the frequency of the input signal for which the peak detector can accurately respond.

One solution, obviously, would be to get a faster op-amp (the Texas Instruments THS3491, for example, has $8000\text{V}/\mu\text{s}$ slew rate, 900MHz bandwidth, 500mA output current, and costs $\$11.80$ each, vs. $\$0.80$ for two op-amps in the TL082) — this may be the only solution if your circuit must respond to very narrow, infrequent pulses. If, on the other hand, you need to track the amplitude of a sinusoid, a simple modification to the design of the peak detector circuit can greatly improve its frequency response (Figure 3-19).

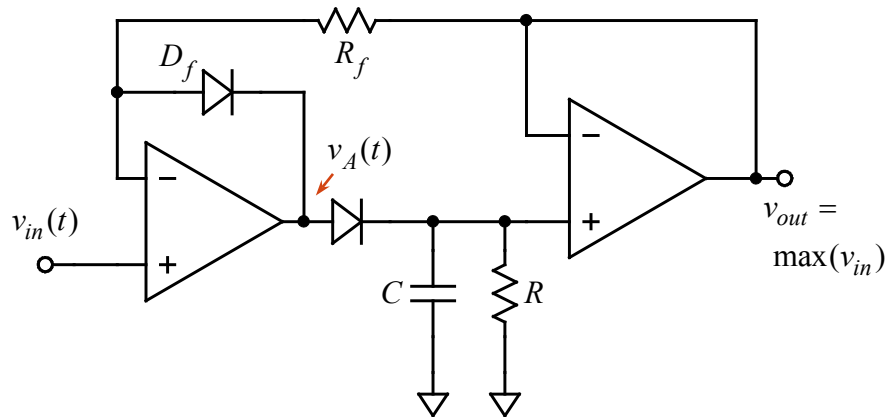


Figure 3-19: A modified peak detector with improved frequency response. The feedback loop now comes from the final voltage follower output; the addition of R_f and D_f ensures that the loop stays closed when v_{in} drops below v_{out} . Now the first op-amp output (v_A) never gets further than a diode drop away from v_{in} while it is below v_{out} rather than saturating at its negative limit as is the case in the Figure 3-18 design.

This modified design adds resistor R_f and diode D_f . First consider the state when v_{in} is rising through the voltage stored by capacitor C so that the original output diode is conducting and C is being charged. Because the voltage follower's output v_{out} matches v_C , then so does the feedback voltage to the first op-amp through R_f . This works because the first op-amp's output voltage v_A is a diode drop higher than v_C when the output diode conducts. Therefore $v_A > v_{out}$ so that diode D_f is reverse-biased. No current then flows through it or R_f , so the first op-amp's $-Input$ voltage, v_- , equals v_{out} , which it works to keep equal to v_{in} (whew!).

Once v_{in} passes its new peak and starts to drop, then $v_{in} < v_{out}$, so the first op-amp's output v_A rapidly decreases, turning off the output diode. Once v_A has fallen a diode drop below v_{out} diode D_f turns on, and the first op-amp's v_- will decrease until it equals v_{in} . The first op-amp

will now maintain this condition by keeping its output v_A a diode drop below v_{in} . The resistor R_f isolates the output voltage v_{out} from the voltage v_- . The current through it is drawn from the second op-amp's output, flows also through D_f , and then into the first op-amp's output.

So whenever the circuit's input voltage is less than the stored peak voltage, the components R_f and diode D_f will keep the first op-amp's output only a diode drop below v_{in} , rather than at negative saturation. When the input rises above the stored voltage, the op-amp's output then needs to slew up through only two diode drops (about a volt) to turn on the output diode and start charging the capacitor. The TL082's output can slew this far in only about $0.1\mu\text{s}$, so its full bandwidth will be available for tracking the input's peak voltage.

A good value for R_f is probably about 22 k or so.

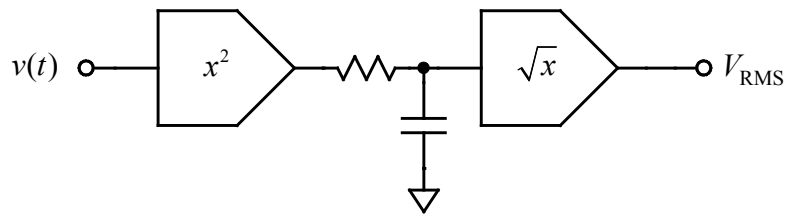
True RMS measurement using analog multipliers

The usual way to determine the power transmitted by an arbitrary, time-varying signal is to measure the average of its squared amplitude (where the average is over a time interval usually much longer than the period of the lowest frequency AC component in the signal). Consequently, the root-mean-square (*RMS*) amplitude of a waveform is typically used to characterize its magnitude:

$$3.4 \quad V_{\text{RMS}} \equiv \sqrt{\overline{v(t)^2}}$$

For a DC voltage, V_{RMS} is just the DC voltage value; for a sinusoid, it is $1/\sqrt{2}$ of the phasor magnitude, or about 35% of its peak-peak voltage. For a complicated signal composed of a DC and several AC components, V_{RMS} is the Pythagorean sum (square root of the sum of the squares) of the individual component *RMS* values.

Since the mean of a time-varying signal is just its DC component (see equation 2.2, page 2-4 of Experiment 2), we can extract the mean of a signal using a low-pass filter with a cutoff frequency well below the lowest frequency of any AC signal component (the effective averaging time, for example, of a simple *RC* low-pass filter is approximately equal to its *RC* time constant). An analog circuit to output the RMS value of a time-varying input signal could then be constructed thusly:



Experiment 3: More circuit ideas

The scale factors of the multipliers used for the square and square root circuits don't matter as long as they are the same; Figure 3-14 on page 3-12 shows a square root circuit using the MPY634.

If you experiment with this circuit, use a variety of input waveforms with various frequencies and RMS amplitudes; compare the circuit's DC output voltage to the signal generator's RMS amplitude setting (make sure that the signal generator output setup is such that it reports the amplitude assuming a "High-Z" load). Estimate the circuit's accuracy and its high frequency response limit.

ADDITIONAL INFORMATION ABOUT THE TEXT IDEAS AND CIRCUITS

Zener diode regulator

Some diodes are designed to be used in their *reverse-bias breakdown region*: Zener and *avalanche* diodes. These types of diodes are very useful as voltage references, simple voltage regulators, and overvoltage protection devices. The Zener and avalanche breakdown effects are described very briefly in the section

Figure 3-20 shows a typical Zener reverse-bias *I-V characteristic curve* — a plot of the relationship between applied reverse-bias voltage and resulting current flow through the diode. The very steep portion of the curve corresponds to the diode's reverse-bias breakdown region; because the curve in this region is so steep, you can see that changes in the diode reverse current correspond to very small changes in reverse-bias voltage. Thus, *the voltage across the Zener diode in this breakdown region is very insensitive to changes in the current through it*. This characteristic makes the Zener diode useful as a simple *voltage regulator*.

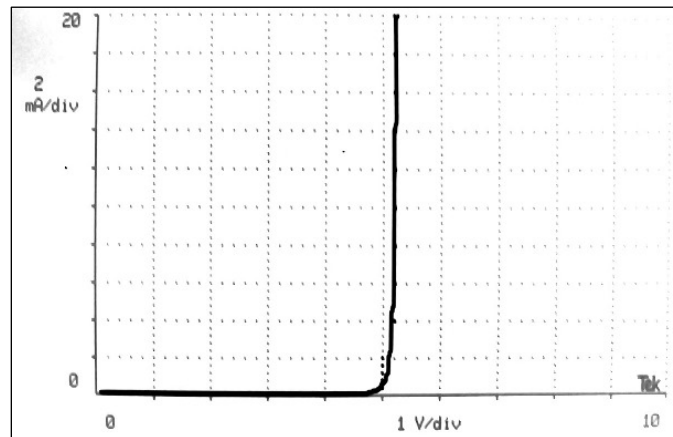


Figure 3-20: Measured Zener diode I-V curve. The reverse diode current is plotted as a function of the applied reverse-bias voltage. As the applied voltage exceeds 5V, the diode current dramatically increases as the diode suffers reverse breakdown. As can be seen from the plot, diode reverse currents above about 15 mA correspond to a reverse-bias voltage of 5.3V. A lab instrument called a *curve tracer* was used to perform this measurement (Tektronix 571).

Assume we have a power supply whose output voltage is greater than the maximum voltage allowed by some device you need to power, and this device needs a *well-regulated power supply voltage* (the voltage is largely unaffected by changes in load current or the input voltage supplied to it). By using a Zener diode whose reverse breakdown voltage is the same as the voltage you need to power your device, you can build a simple voltage divider circuit which will satisfy the requirements for your device's power supply.

Consider the circuit in Figure 3-21, which is just a voltage divider with a reverse-biased Zener diode as the bottom element (note how the cathode end of the schematic symbol for a Zener diode differs from that for a regular diode). If the input voltage exceeds the Zener diode's breakdown voltage, then the diode may break down, with voltage V_R across it. Thus $v_{out} = V_R$. Now the voltage across the resistor R is determined: $v_{in} - V_R$, and the current through the resistor is also known: $i_{in} = (v_{in} - V_R)/R$. This current is divided between the current through the Zener diode, i_Z , and that through the load, i_{load} . Thus $i_Z = i_{in} - i_{load}$.

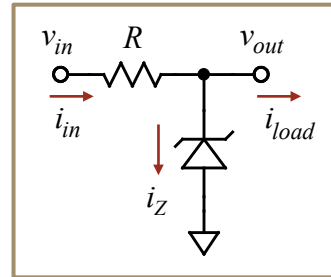


Figure 3-21: Zener diode voltage regulator circuit. The input voltage v_{in} must be greater than the Zener diode's reverse breakdown voltage, V_R . This large voltage causes the diode to suffer reverse breakdown, with the difference $v_{in} - V_R$ imposed across resistor R , establishing the input current i_{in} . This current is divided between that required by the load, i_{load} , and the reverse current through the diode, i_Z .

Note how this circuit holds the voltage applied to the load constant (equal to V_R) even if i_{load} or v_{in} varies. Changing i_{load} doesn't change v_{out} or i_{in} because the voltages across the diode and across R don't change; so for changes in i_{load} , $\Delta i_Z = -\Delta i_{load}$ (we assume here that the Zener's I-V characteristic curve is very steep in its breakdown region: $dV_R/di_Z \approx 0$). Similarly, changes in v_{in} don't affect v_{out} because, although i_{in} changes, $\Delta i_Z = \Delta i_{in} = \Delta v_{in}/R$, but the steep Zener I-V characteristic ensures that v_{out} is nearly unaffected by this change in i_Z .

DESIGNING A ZENER VOLTAGE REGULATOR CIRCUIT

This example will illustrate how you would design a voltage regulator using the Zener diode circuit in Figure 3-21. Assume the load is a digital circuit which requires a stable voltage of no more than 5.5V in order to operate (the so-called CMOS *digital logic family* to be discussed in a later experiment would fall into this category), and you wish to power it from a 9V battery. The circuit will require a minimum of 10mA, but could draw as much as 20mA when an LED (light emitting diode), a part of the circuit, is illuminated. You want the circuit to work properly even if the battery has discharged to the point where it can only supply 7V.

Here are the basic Zener regulator design steps:

1. Choose a Zener diode breakdown voltage. The Zener diode whose characteristic curve is shown in Figure 3-20 should work, since its nominal breakdown voltage of 5.3V is near, but less than, the specified 5.5V limit.
2. Choose a *minimum* Zener current which will give good voltage regulation (steep I-V curve); according to Figure 3-20, a current of 10mA is well into the diode's breakdown region, so we will design the circuit for that minimum diode current.

3. The minimum Zener current + the maximum load current = the design target for the current i_{in} through the resistor R (Figure 3-21). Therefore, for this example, design $i_{in} = 10\text{mA} + 20\text{mA} = 30\text{mA}$. The resistor value is then chosen to give the correct voltage drop at this design current when the input voltage is at its *minimum* (7V for this example). Thus $R = (7\text{V} - 5.3\text{V})/30\text{mA} = 56.7\Omega$. Choose the closest standard resistor value, which, in this case, is $R = 56\Omega$.
4. Now consider what would happen when the circuit experiences the opposite extreme: maximum source voltage v_{in} (9V) along with minimum load current i_{load} (10mA). Since the Zener diode will keep the output voltage at 5.3V, the voltage drop across R is now $9\text{V} - 5.3\text{V} = 3.7\text{V}$, and with the chosen value for R , $i_{in} = 3.7\text{V}/56\Omega = 66\text{mA}$. Under these conditions the diode current will have increased to 56mA ($i_z = i_{in} - i_{load}$).
5. Use the results from (4) to determine the worst-case power dissipations in the resistor and the diode:

$$\text{Resistor: } P = VI = 3.7\text{V} \times 66\text{mA} = 0.24\text{W}$$

$$\text{Diode: } P = VI = 5.3\text{V} \times 56\text{mA} = 0.3\text{W}$$

These results specify the required minimum power dissipation capabilities of the components, which should be at least 150% of the calculated values (to be on the safe side).

Note that the power required from the 9V battery is $9\text{V} \times 66\text{mA} = 0.6\text{W}$, while the load may be consuming only $5.3\text{V} \times 10\text{mA} = 0.05\text{W}$, so the efficiency of our simple regulator is a measly $0.05/0.6 \approx 8\%$ (ouch!). This is typical for a Zener voltage regulator, which is called a *shunt regulator*: the current drawn from the source is always greater than the maximum required load current, even when the load is drawing little or no current most of the time.

Zener diode voltage regulator circuits are really only practical when both the variation in the required load current (i_{load}) and the expected variation in the source voltage (v_{in}) are small.

For this particular problem a more efficient solution would be to use a *series regulator* (the standard type implemented by special-purpose voltage regulator ICs), or, even better, a *DC-DC converter*, which is a type of *switching regulator* that can convert 9V power to 5V power with efficiencies exceeding 80%.

LEDs

A *light emitting diode* (LED) is a PN junction diode made from a semiconductor material with a relatively large energy gap. Consequently, when an electron and hole recombine the energy released may be carried away by a visible light (or near infrared) photon. Otherwise, LEDs behave much the same as any other semiconductor diode. Because of its larger gap

voltage, the forward-bias voltage drop for an LED is significantly higher than for a silicon diode; typical forward voltages for various LEDs are:

Table 3-1
Typical LED Forward Voltage Drop (10mA forward current)

IR (800-1000nm)	Red (620-680nm)	Yellow (570-590nm)	Green (520-565nm)	Blue (450-470nm)
1.2V	1.8V	2.0V	2.2V	4.0V

White LEDs are usually constructed using a blue emitter with a fluorescent coating, so their forward voltage drop is the same as for blue LEDs.

Caution

LED reverse breakdown voltage can be only 5V! Be careful to avoid using an LED in a circuit which could cause it to suffer reverse breakdown.

Since the LED is a PN junction diode, *its forward current must be limited by the external circuit*; this is most often accomplished using a resistor in series with the diode.

The intensity of the light output of an LED is very nearly proportional to its forward-bias current; 10 – 15 mA is usually more than sufficient to provide a nice, bright indicator light. Be careful to avoid excessive forward currents or the LED will quickly fail; if a very bright source is needed, specially designed and cooled LEDs are available (at a price). LED indicators are often controlled by digital logic circuitry, and are turned off or on depending on the logic circuit state; we'll discuss these sorts of circuits in later experiments. In Experiment 1 we described an op-amp based *transconductance amplifier* which supplied an LED with current proportional to a supplied input voltage, thus avoiding the nonlinearity introduced by the diode's current-voltage relationship.

Often you will want an LED *pilot light* to illuminate whenever the power supply to a circuit is activated (like the LEDs on the circuit trainer breadboard). Some circuits to do this job are shown in Figure 3-22. For circuits (a) and (b) from that figure, the minimum voltage required for any significant illumination is the LED's turn-on voltage (Table 3-1). The resistor value is chosen so that the desired current will flow through the diode to give the desired intensity. For example, illuminating a green LED with 10 mA from a 12V source would require a resistor value of $R = (12V - 2.2V)/10mA \approx 1k\Omega$; a 5V source would require $R \approx 270\Omega$.

The circuit in (c) is interesting because it uses a Zener diode to set the minimum voltage required for the LED to illuminate. This would be useful if you don't want the pilot light to illuminate if the power supply voltage is less than some threshold value. For example, assume you are using a 9V battery for your power supply, but the circuit needs at least 7V to

operate properly. If the Zener breakdown voltage is 5.1V and you use a red LED (turn-on 1.8V), then at least $5.1\text{V} + 1.8\text{V} = 6.9\text{V}$ would be needed to illuminate the LED. With a current of 10mA at 9V, $R = (9\text{V} - 6.9\text{V})/10\text{mA} \approx 200\Omega$.

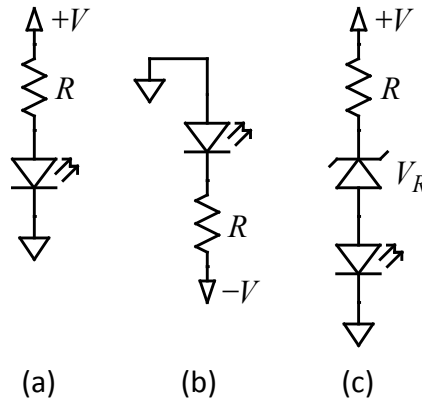


Figure 3-22: LED pilot light circuits. (a) positive supply voltage; (b) negative supply voltage; (c) positive voltage, but with a minimum voltage threshold for LED illumination using a Zener diode. See the text for details.

In experiment 4 we'll look at a more precise method for illuminating an LED only if some voltage threshold is crossed.

Using a PN junction for temperature sensing

The PN junction's temperature sensitivity (equation 3.10 on page 3-42) makes the semiconductor diode a useful temperature sensor; in this section we present one example of a temperature monitor circuit. If a silicon diode is forward-biased so that its current is much greater than I_R in equation 3.10, then the I-V relationship is:

$$3.5 \quad I \approx I_0 \exp \left[\frac{-q_e \left(\frac{V_g - V}{T} \right)}{\eta k_B} \right]$$

This implies that if a diode forward-bias current is held constant at, say, 0.1mA ($\sim 10^5 I_R$), then $V_g - V \propto T$. For a silicon diode at room temperature, $(V_g - V)/T \approx 2\text{mV/K}$, and this is the sensitivity of the forward-bias voltage to a change in junction temperature. Here is a circuit:

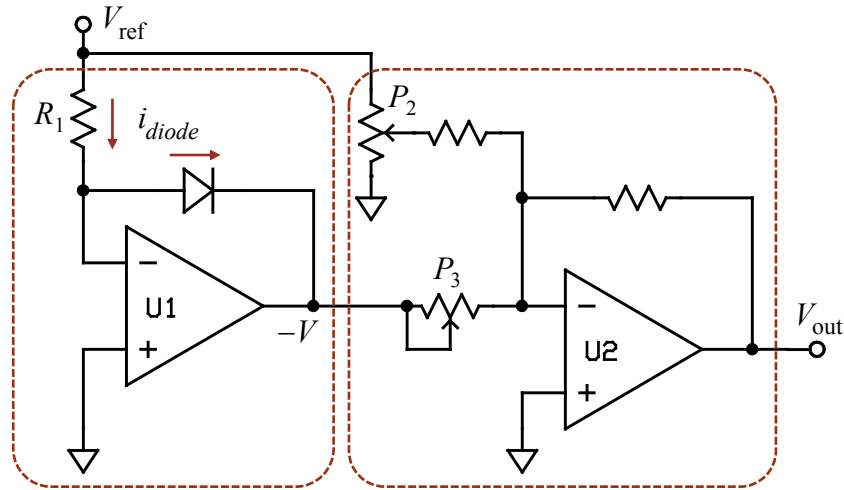


Figure 3-23: Temperature circuit using a diode sensor. Amplifier U1 maintains a constant diode forward-bias current of V_{ref}/R_1 ; the amplifier output is then $-V$, where V is the diode's forward-bias voltage. The inverting, summing amplifier U2 scales and offsets this voltage to provide the circuit's output (potentiometer P_2 trims the output voltage's DC offset, P_3 trims $\Delta V_{out} / \Delta T$).

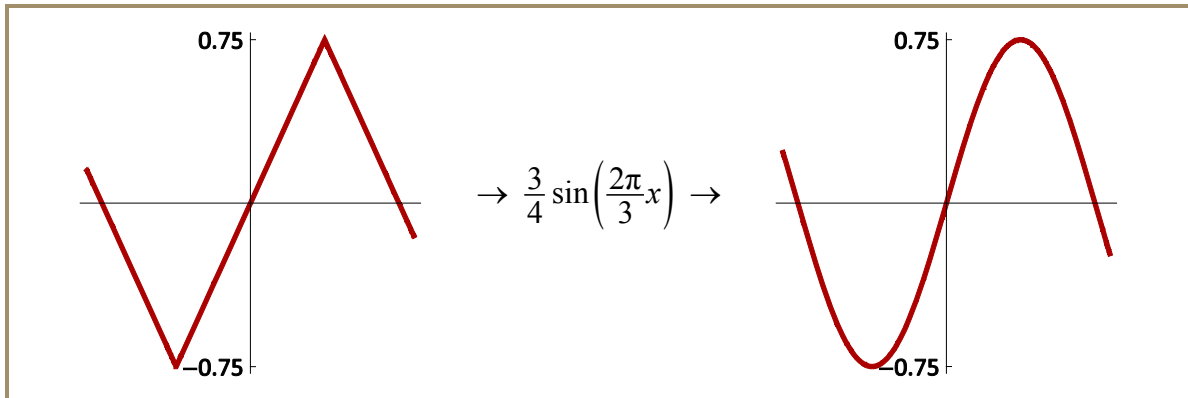
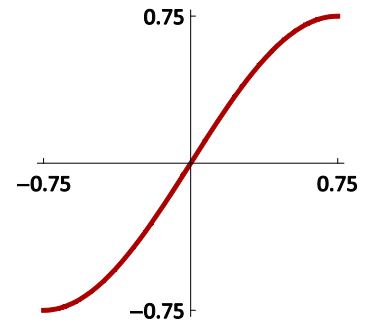
To provide an output temperature response of $0.1V/^\circ C$, the gain of the U2 inverting amplifier stage should be ≈ 50 when R_3 is set to its center position; the summing amplifier's gain for the offset adjust input (R_2) should be 1 or less, so when R_2 is set near its center position the output corresponds to the diode temperature. There are several ways this circuit could be refined, so that constructing and calibrating such a circuit could be a good final project. The [Texas Instruments LM35](#) series of analog temperature sensor ICs implements a version of the circuit in Figure 3-23 in a single, small, calibrated device.

Approximating a transcendental function using analog multipliers

In Experiment 4 we'll see how to use a couple of op-amps and a few other components to construct a "function generator" which can output accurate square and triangle waveforms and whose frequency is easy to adjust with a single potentiometer. Unfortunately, constructing a *sine-wave* generator whose frequency is also easy to adjust is not so simple. Therefore the design of an analog function generator which can output a sinusoid as well as square and triangle waveforms offers an interesting challenge.³ The solution to this design problem is to use a relatively simple nonlinear circuit which can distort a basic function generator's triangle-wave output into a shape which offers a good approximation of the desired sinusoid. In this section we show one way this can be achieved using a pair of multipliers.

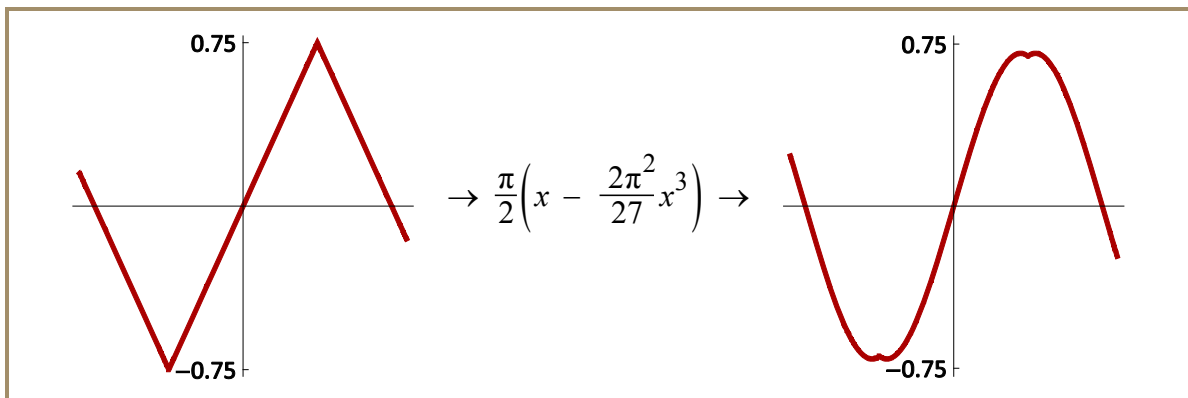
³ By "analog function generator" we mean one that is constructed from analog integrated circuits such as op-amps and analog multipliers, as opposed to digital circuitry with special-purpose, computer-like subsystems. The lab signal generators use this latter approach and accomplish their tasks using a process called *direct digital synthesis*, a relatively modern technique which has all but completely replaced analog designs.

Consider calculating an approximation to the transcendental function $(3/4)\sin(2\pi x/3)$ for $-3/4 \leq x \leq 3/4$, proportional to x near 0 but flattening out as $|x|$ approaches $3/4$, as shown in the figure at right. This particular domain and range will be ideal for a triangle-wave to sine-wave converter, because by applying this function to a triangle wave as it oscillates linearly between $-3/4$ and $+3/4$, the converter circuit's output will produce a sinusoid over the same range as shown below (why we chose $3/4$ as a limit rather than 1 will become apparent later).



At first you might be tempted to design a circuit to calculate a couple of terms of a traditional Taylor series expansion of the function about $x = 0$, but it turns out that one can often do better with a more clever choice for the approximating function. We begin with the two-term Taylor expansion:

$$\frac{3}{4}\sin\left(\frac{2\pi}{3}x\right) = \frac{\pi}{2}\left(x - \frac{2\pi^2}{27}x^3\right) + O[x^5]$$

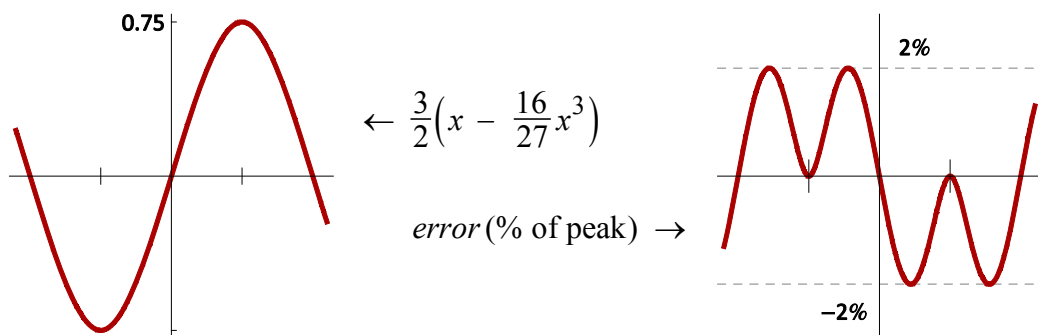


The problem with using this function should be obvious from the above plot. Although the Taylor series converges quite rapidly near $x = 0$, its convergence at $|x| = 1$ is slow.

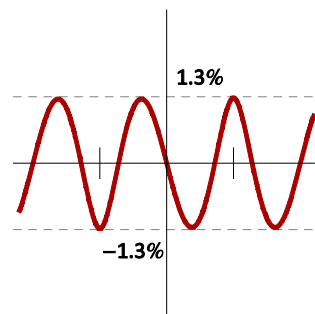
Maybe a slightly different cubic approximation will fare better. How about

3.6
$$\frac{3}{4} \sin\left(\frac{2\pi}{3}x\right) \approx \frac{3}{2}\left(x - \frac{16}{27}x^3\right); \quad -\frac{3}{4} \leq x \leq \frac{3}{4}$$

This function matches the sine at its peaks as well as at its zero crossings and it also has vanishing slope at $|x| = 3/4$, as does the sine function. Note also that the coefficients have rational values which actually aren't very different from those in the truncated Taylor expansion. Here's how it fares when operating on the triangle waveform:



So the average error of this function is only about 1% of its peak amplitude with a maximum error of 2%. It really looks like a sinusoid as well with no obvious artifacts. The function is actually even better than the above error would indicate. When compared to a sinusoid with a very slightly lower amplitude (decreased by less than 1.3%), its maximum error is cut nearly in half, as shown at right. By not exactly matching that sine-wave's amplitude, the function better matches it at intermediate positions of the waveform. Because a typical analog multiplier has an error of 2%, it is a waste of time to try to do any better than this.



Now to implement (3.6) using analog multipliers. Both the MPY634 and the AD633 have scale factors of 10 V. This means that not only does a voltage of $v_x = 10\text{V}$ correspond to $x = 1$ in our equations, but also that 10V is the maximum allowable input and output voltage magnitude for accurate multiplier operation. This explains why we chose our model to require that $|x|$ only go to $3/4$, corresponding to $v_x = 7.5\text{V}$: we have some elbow room to

make minor amplitude tweaks to the input triangle wave without hitting the multipliers' 10V limitation.

The circuit is shown in Figure 3-24. The first multiplier, shown in its simplified, generic form, is configured as a simple squarer, taking $x(t)$ and generating $x^2(t)$. Its output goes to a voltage divider whose output ratio $a = 16/27$, the coefficient of the cubic term in (3.6). One possible choice of standard resistor values to achieve this could be 11k and 16k.

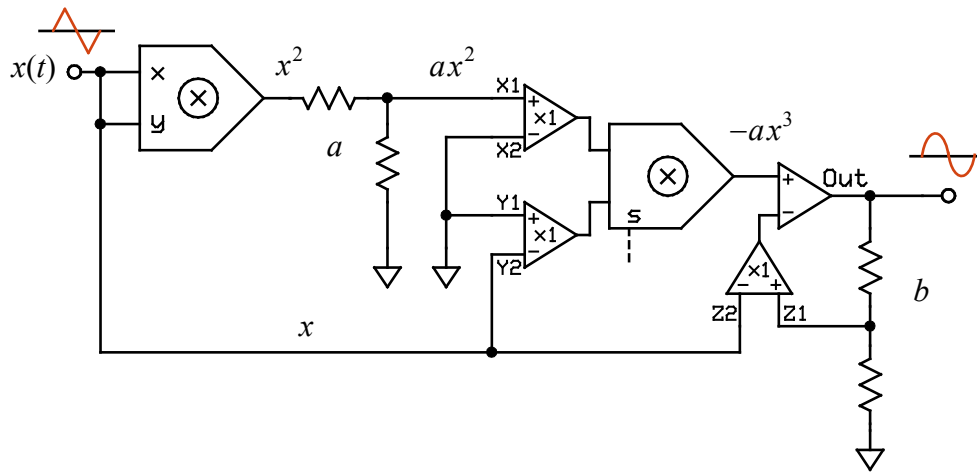


Figure 3-24: The triangle-wave to sine-wave converter circuit. Its operation is described in the text.

The second multiplier in Figure 3-24 is shown as the MPY634 in all of its detail. We're going to make good use of its complicated circuitry to complete the calculation in (3.6). It takes ax^2 from the first multiplier and multiplies by $-x$ to complete the calculation of the cubic term in (3.6). The IC's Z2 input is then used to add x to this result. The feedback network around the MPY634's output op-amp is chosen to perform the final scaling operation of multiplying by $b = 3/2$. The lower resistor in its feedback voltage divider should have a value twice that of the upper, so 10k and 20k might be good choices.

Finally, the circuit expects a triangle-wave input with an amplitude of 7.5 V (15 V peak-peak), which corresponds to 0.75 for the peak amplitude of x in (3.6). Because of resistor tolerances in the voltage divider a and tolerances in the multipliers' accuracies, the triangle wave amplitude (and possibly DC offset) should be adjusted to achieve the best result. An appropriate method to evaluate the resulting output's deviation from a sinusoid is to use the FFT mode of an oscilloscope and then adjust the input to minimize the higher harmonics evident in the output's FFT.

Full-wave and bridge rectifiers

When used to rectify an AC power source to generate a DC power supply voltage, the simple, half-wave rectifier of Figure 3-7 on page 3-6 suffers from an obvious flaw: its single

diode discards half of each AC waveform cycle as the circuit attempts to keep power supplied to its load. To keep the residual AC *ripple voltage* in the circuit's output reasonably small, the filter capacitor used to maintain the circuit's output must be large so that its RC time constant is much longer than an AC cycle. Examine Figure 3-7 again: if the decay in the capacitor's voltage must be small, then the time during which the diode conducts as the source voltage rises though the capacitor's voltage will be a very small fraction of an AC cycle. Consequently, a very large current must flow from the source through the diode in order to recharge the capacitor, because the average current from the source must equal that supplied to the load in order to maintain a steady output.

Using both the negative and positive peaks of the input AC waveform could help mitigate this problem by roughly halving the time between successive capacitor recharges. In other words, a full-wave rectifier would be more desirable than a simple half-wave model. Its implementation is shown in Figure 3-25. A *power transformer* is used to convert the incoming AC voltage to a value suitable for the desired DC output voltage (transformers are

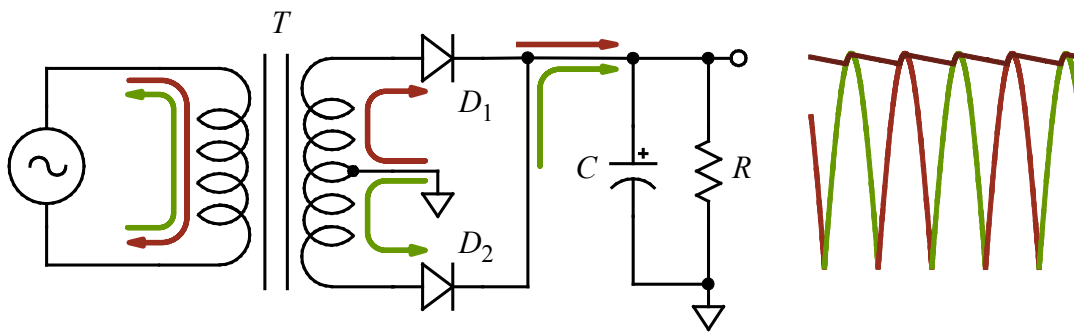


Figure 3-25: Transformer-coupled full-wave rectifier. The AC source drives transformer T which has a center-tapped secondary coil. As the driving current reverses direction on alternate half-cycles, so does the induced voltage on the transformer's secondary. This voltage in turn alternately forward biases one or the other of the two diodes. Thus current always flows in the same direction through the load R , returning to the transformer secondary's center tap via the output circuit ground connections. Unfiltered (red and green) and filtered (brown) output voltage waveforms are shown at right. The maximum output voltage across R will be $1/2$ of the total secondary coil peak voltage (minus a diode drop).

described in the *Additional Information* section of the Experiment 2 text). The oscillating current in the transformer's primary coil induces an output voltage in its secondary coil. referring to Figure 3-25, we see that the transformer T 's secondary coil is *center-tapped*, so that equal voltages are induced across the two halves of the coil, with the same polarity between the coil's bottom end and its center tap and between the center tap and the coil's top. As shown in the figure, the secondary's center tap will become the output circuit's ground return, whereas the two ends of the coil supply power through their respective diodes to keep the capacitor C charged.

To better understand the operation of the full-wave rectifier circuit, first consider the case where the capacitor C in Figure 3-25 is not installed (the *unfiltered* case). The transformer secondary's center tap is connected to each end of the coil through the load resistor R and the two diodes D_1 and D_2 . During each half cycle of the AC source, the voltage across the secondary will forward bias only one of the two diodes while the other is reverse-biased. Current will flow through the forward-biased diode, through resistor R , and return to the coil through its center tap. On the next half cycle, current will flow through the other diode. In either case, current will flow into the top of resistor R and back into the center tap as shown by the red and green arrows in Figure 3-25. The peak output voltage across R will be the induced voltage *across only one half of the secondary*, minus a diode's forward voltage drop. Adding the capacitor C will smooth the output voltage across R in the same way as for a half-wave rectifier, but now the capacitor's discharge time between recharges will be roughly halved.

An op-amp power supply such as that used in the circuit breadboards requires a *bipolar* DC power supply such as $\pm 12\text{V}$. This can be achieved with the full-wave circuit by simply adding another output circuit to the transformer's secondary in parallel with the original circuit. Reversing the diode orientations in this parallel circuit will generate the opposite-polarity voltage as shown in Figure 3-26. This is exactly the sort of circuit used in the breadboard power supply. The schematic shows that the lower left diode is connected to the upper end of the transformer's secondary coil – the little “jump” in the schematic's vertical connecting wire indicates that it makes no connection to the horizontal wire coming from the secondary's lower end.

Now consider what would happen if the load attached to the dual output circuit is connected between its two output terminals with no connection to the transformer secondary's center tap. The total voltage across the load would then be the secondary coil's total voltage minus two diode drops (one from each conducting output diode). In this case the center tap would

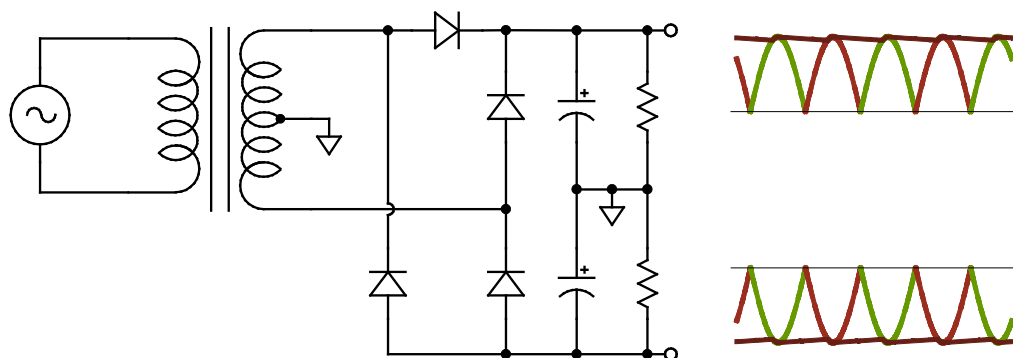


Figure 3-26: Dual-output full-wave rectifier. The center-tapped secondary feeds two separate full-wave rectifier output circuits in parallel. The diodes of the lower circuit are oriented oppositely to those in the upper circuit, so the lower circuit's output will have negative polarity with respect to the circuit ground. The unfiltered and filtered outputs are shown with colors corresponding to those in Figure 3-25.

be unnecessary; the two series-connected filter capacitors could be replaced by a single capacitor and the two resistors by a single load resistor. The result is the *bridge rectifier* circuit, shown in Figure 3-27.

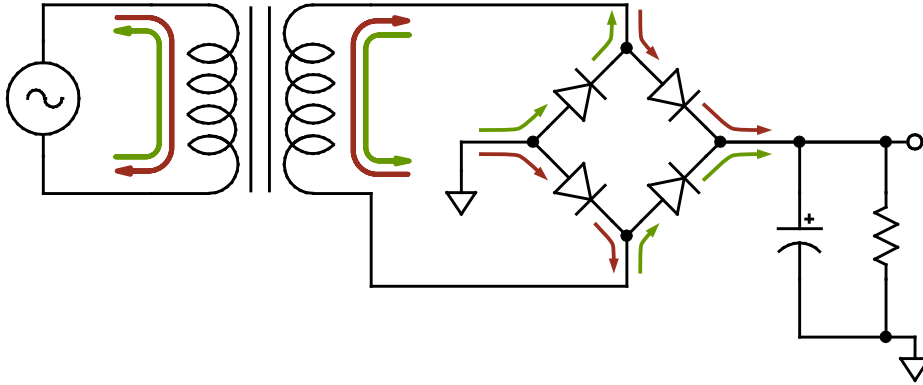


Figure 3-27: Transformer-coupled bridge rectifier. This circuit is derived from the dual-output full-wave rectifier shown in Figure 3-26, but it has only a single load resistor and filter capacitor connected between that circuit's two opposite-polarity output voltages. As the driving current reverses direction on alternate half-cycles, so does the induced voltage on the transformer's secondary. This voltage in turn alternately forward biases one or the other of two diode pairs in the four-diode "bridge," as shown by the red and green sets of arrows. The maximum output voltage across the load will be the total secondary coil peak voltage (minus two diode drops).

The four-diode circuit topology found in both Figure 3-26 and Figure 3-27 is called a "diode bridge" and is conventionally drawn using the diamond-shaped arrangement of Figure 3-27. Single packages containing four diodes are readily available and are called *bridge rectifiers*.

THE PHYSICS OF THE PN JUNCTION DIODE

Insulators, conductors, and semiconductors

The electrical conductivities of solid materials for the most part fall into one of two classes: conductors and insulators (metals make up most of the conductors, and nonmetals are usually insulators). Although all materials are very nearly electrically neutral (equal numbers of protons and electrons, so that they carry no net charge), the nature of the chemical bonds which bind the atoms or molecules of a solid to one another determines its class of electrical conductivity.

An atom's *valence electrons* — the outer, most weakly bound electrons — are the ones which participate in chemical bonding. The atomic nucleus along with the much more strongly bound inner electrons comprise a positively-charged *ion core* which remains intact and is surrounded by the interacting valence electrons. The chemical bonding process causes these many ion cores in a solid to arrange themselves in a mostly regular, crystalline structure. This regular, periodic array of positively-charged cores creates a similarly regular, periodic electrostatic field within which the myriad valence electrons move.

The quantum-mechanical nature of these microscopic, negatively-charged particles (the valence electrons) as they evolve in the periodic electrostatic potential of the ion cores requires that they each occupy a state of motion (and total energy) in one of several distinct *energy bands*, analogous to the quantized energy states an electron may occupy in a single atom or molecule. The width of a typical energy band is on the order of a few to several electron volts (same order of magnitude as the binding energy of a valence electron in one of the atoms), and adjacent energy bands are often separated by a similar energy, although they may also overlap. Each band has enough distinct quantum states to contain twice the number of electrons as there are molecules in the macroscopic solid crystal (i.e. $\sim 10^{23}$).

Room temperature ($\approx 290\text{K}$) corresponds to random, thermal particle energies of $\sim 1/40\text{eV}$ (electron volt), much smaller than the width of an energy band but much larger than the energy spacing between the individual states in a band. Because electrons are subject to *Pauli Exclusion* (each electron must be in a unique, distinct quantum state), the valence electrons of all the various atoms in a solid fill the available states starting with the lowest available energy. Because room temperature corresponds to a fairly small energy, the energy of the topmost filled states is fairly well-defined and is called the electrons' *Fermi energy*. Random thermal jostling can only affect the states of individual electrons with energies near the Fermi energy, because those with much lower energies are surrounded by quantum states already occupied by other electrons, so they're stuck in their current states.

Now, one of two situations can occur for our valence electrons in a solid:

- (1) The number of electrons is such that they exactly fill all the available states in some number of energy bands, and higher energy bands are completely empty.
- (2) One energy band (or possibly more, if some bands overlap) is only partially filled and has many unoccupied states still available; all other bands are either completely filled or completely empty.

Electrons occupying a completely filled energy band do not participate in electrical conduction. The reason for this is that such a band corresponds to all physically possible states of individual electron motion in all directions consistent with the energies of the electrons in that band. Applying an external electric field doesn't change this situation unless the field is so intense that it can cause electrons to transition to another (partially filled or empty) energy band. Thus, no new net motion of electrons can be induced by the presence of the field, so the electrical conductivity contributed by a completely full (or, of course, completely empty) energy band is zero.

This last result implies that solids with situation (1) above are *insulators* (or maybe semiconductors). Since each band has twice the number of states as there are molecules in the crystal, insulating materials most often arise when there is an even number of valence electrons participating in the chemical bonding forming the solid. Situation (2), on the other hand, allows electrical conduction to proceed using the electrons in the partially-filled energy band. Electrons near the Fermi energy in the band have a wide selection of nearby empty states, so an applied electric field can accelerate them, and their resulting motions can carry a net flow of charge (electric current) through the solid. These materials are *conductors*, and partially-filled energy bands are characteristic of the so-called *metallic bond*.

Semiconductors have valence electrons whose situation falls into category (1): bands containing electrons are completely filled, at least at cold temperatures. What makes them different from insulators, however, is that the bottom of the nearest empty energy band (called the *conduction band*) is only about an eV or so away from the top of the highest-energy filled band (the *valence band*). Consequently, random thermal jostling of the ions in the lattice can impart enough energy to a few electrons with energies near the top of the valence band to excite them into levels near the bottom of the conduction band. In this case both the valence band and the conduction band become *partially* occupied (although just barely), and the material becomes a poor conductor (poor because only a tiny fraction of the valence electrons get bumped up into the conduction band). The higher the temperature, the greater the number of valence electrons thermally excited into the conduction band — the number goes as:

3.7

$$n_i \propto T^{3/2} e^{-E_g/(2k_B T)}$$

where E_g is the magnitude of the energy gap between the valence and conduction bands and k_B is Boltzmann's constant. In the case of silicon, this amounts to $\sim 10^9$ electrons per cm^3 at room temperature (compare with copper's 0.8×10^{23} per cm^3).

The archetypal semiconductors are the elements silicon ($E_g = 1.12 \text{ eV}$) and germanium ($E_g = 0.67 \text{ eV}$), each of which forms a diamond lattice with four covalent bonds per atom. Carbon in its diamond form ($E_g = 5.5 \text{ eV}$) is beginning to find applications in solid-state devices, but its large energy gap makes it more properly classified as an insulator.⁴ Several compounds and alloys form commercially important semiconductors, including GaAs (the first LED, emitting in the near infrared), InP, GaAsP, and InGaN (blue LED).

Electrons and holes; impurities and doping

The diagram at right illustrates the distribution of electrons between the top of the valence band and the bottom of the conduction band for a semiconductor at a fairly high temperature (so that there have been a considerable number of electrons excited into the conduction band). The density of quantum states grows as $\sqrt{\Delta E}$ as you move away from the band edges, as shown by the right-hand curves in the figure. The Boltzmann factor $\exp(-\Delta E/k_B T)$ gives the relative probability that any one state is occupied in the conduction band or unoccupied in the valence band. Because the number densities of the conduction electrons and the holes in their respective bands are low (much smaller than the crystal's atomic number density), the charge carriers will distribute themselves in accordance with the classical Maxwell distribution, so the overall occupation densities go as the left-hand curves in the figure.

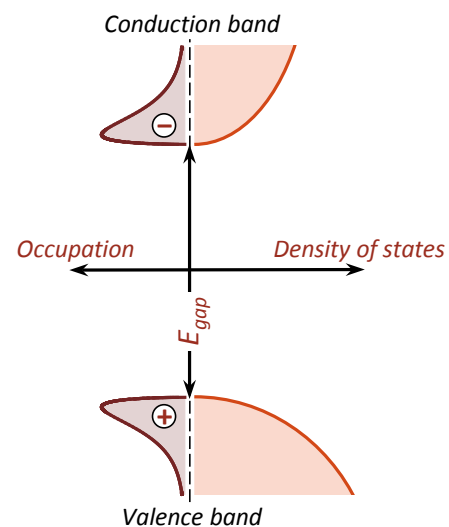


Figure 3-28: Densities of states and occupations by electrons (conduction band) and holes (valence band) for a pure semiconductor (only intrinsic charge carriers); energy increases in the vertical direction. The kinetic energy distribution of the charge carriers in each band is classical.

The dynamics of the relatively small number of electrons in the conduction band is very well approximated by treating them as classical particles (with a negative charge of $-q_e$, of course), but their *effective mass* is determined by the shape of the density of states curve near the bottom of the conduction band: the sharper the curve near the minimum, the lighter the effective mass. In the valence band, only a small fraction of the states near its top are

⁴ Carbon in the form of graphite has layers of honeycomb arrangements of atoms (a single layer forms a 2-D structure called *graphene*). Graphite is a *semi-metal*: it has partially-filled bands like a conductor, but these bands are very nearly empty. Consequently the conduction electron density in graphite is much smaller than that of a typical conductor, 10^{18} vs. 10^{23} electrons per cm^3 .

unoccupied. Interestingly, the dynamics of the remaining electrons near the top of the valence band are such that they have a *negative effective mass*, since the density of states *decreases* with increasing energy near the top of the band.

The consequence of this unusual electron behavior near the top of the valence band is that the unoccupied quantum states evolve as though they were *positively charged particles* ($+q_e$) *with a positive effective mass and with energies increasing as they move further down from the band top!* These “positive charge carriers” near the top of the valence band are called *holes*. Thus, when an electron is excited from the valence band to the conduction band, *two* charge carriers are created: the electron ($-q_e$) and the hole ($+q_e$) it left behind. Since the energy an electron must gain to cross the energy gap between valence and conduction bands is at least E_g , but two “particles” were created by this transition, *the required energy per particle* is $E_g/2$ — this observation explains the 2 in the Boltzmann factor in equation 3.7.

Thus a pure semiconductor has a conductivity which is a very strong function of temperature, rising rapidly as temperature increases (equation 3.7). This effect is used to make a *thermistor*: a resistor with a large, *negative temperature coefficient* (decreasing resistance as temperature rises) which acts as a very sensitive, fast-acting temperature sensor for the range of about -100°C to $+150^\circ\text{C}$.

Semiconductor materials are custom-made to be much more flexible and useful through the process of *doping*: introducing various amounts of impurity atoms into the semiconductor crystal which have a different valence than the semiconductor. For example, mixing a small amount of phosphorous (valence 5) into a silicon crystal will introduce atoms each with an extra valence electron left over after it forms bonds with surrounding silicon atoms. What would be the consequences of these extra electrons to the physics of the material? It turns out that the energy of this extra valence electron is very close to the energy of the bottom of the conduction band (in the case of P in Si, the energy is only 0.044 eV below the conduction band). If there are relatively few of these *donor* impurity atoms, then it is very likely that such electrons will eventually be thermally excited into the conduction band: once there they move away from the impurity atoms and are unlikely to recombine with them.

So even if the ambient temperature is cool enough that almost no electrons would be excited from the valence band to the conduction band, electrons from donor impurities will nearly all find their way into the conduction band, providing a largely temperature-independent cadre of negative charge carriers ($-q_e$) along with the same number of fixed, positively-charged ions ($+q_e$) distributed throughout the crystal lattice. Such a material is called an *N-type semiconductor*.

Similarly, introducing a valence 3 impurity atom (such as aluminum into silicon) will leave an unsatisfied bond because of the missing electron. Again, the energy required to promote a nearby valence electron into this spot is small compared to the semiconductor’s energy gap (0.057 eV for Al in Si). Thermal agitation will eventually do the trick, and the vacated

valence state becomes a hole which quickly moves away from the impurity atom, trapping the promoted electron at the impurity site. Thus these *acceptor* impurity atoms become fixed, negatively-charged ions ($-q_e$) in the lattice, whereas an equal number of holes form a temperature-independent group of positive charge carriers ($+q_e$), creating a *P-type semiconductor*.

Adding dopants to a semiconductor can not only introduce charge carriers (called *extrinsic* charge carriers), but will also suppress the thermal creation of electron-hole pairs described by equation 3.7 (called *intrinsic* charge carriers). This is because the product of the number of conduction electrons (n_c) and the number of holes (p_v) is related to the number of intrinsic charge carriers thermally created in a pure (undoped) semiconductor (n_i) by the laws of statistical mechanics:

$$3.8 \quad n_c p_v = n_i^2$$

For example, the addition of 10 parts per billion phosphorous to a silicon crystal would introduce 5×10^{14} extrinsic conduction electrons per cm^3 ; with $n_i \sim 10^9$ electrons per cm^3 , we see that there will be only $p_v \sim 2000$ holes per cm^3 ! These holes are called *minority carriers* in the N-type silicon under discussion; the conduction electrons are the *majority carriers*. Since for this example $n_i \ll n_c$, the temperature dependence of n_c will be quite small, so, from equation 3.8, $p_v \propto n_i^2$. Thus *the temperature dependence of the minority carriers is very large*: from equation 3.8,

$$3.9 \quad \langle \text{minority carrier density} \rangle \propto T^3 e^{-E_g/(k_B T)}$$

The equilibrium PN junction

Now consider the case of a semiconductor crystal with inhomogeneous doping. As a concrete (but quite artificial) example, assume that we take a single P-type crystal and a single N-type crystal and then join them along a planar boundary so as to form a single crystal with an abrupt change in doping at this boundary. The result is a *PN junction* (Figure 3-29) at the interface between the two semiconductor types.

Far from the boundary the charge carrier densities must approach their homogeneous, thermal equilibrium values. Near the interface, on the other hand, the large gradients in the hole and conduction electron densities will drive diffusion of these charge carriers across the boundary, where they will eventually recombine with carriers of the opposite sign. The reduced majority carrier densities near the boundary induce a net charge density and resulting electric field near the interface because of the now unbalanced charge of the impurity ions in each semiconductor. This electric field will repel the majority charge carriers on either side of the boundary, and an equilibrium condition is reached preventing further net diffusion of carriers across the boundary (bottom illustration in Figure 3-29). The electric field near the boundary generates a potential difference between the P-type and N-type sides of the junction, with the N-type material at the higher potential. This *contact potential* of the PN junction is very close to the *gap voltage*: $E_g/q_e \equiv V_g$ (1.12 V for silicon). The result, as we shall see, is the creation of the *PN junction diode*.

It turns out that the equilibrium situation will be attained only when a region near the PN interface is almost completely devoid of charge carriers: the so-called *depletion layer*, as shown in Figure 3-29. The charge density in this region is then given by the number densities of the impurity ions on each side of the boundary, which are nearly equal to the corresponding majority carrier number densities far from the boundary. Since the potential changes by $\sim V_g$ across the depletion layer, it is straightforward to calculate its equilibrium width, which will typically be in the range of $10^2 - 10^4 \text{ \AA}$ (about 2000 \AA for silicon with a part per million doping), and the magnitude of the electric field at the interface is in the range of $10^5 - 10^7 \text{ V/m}$.

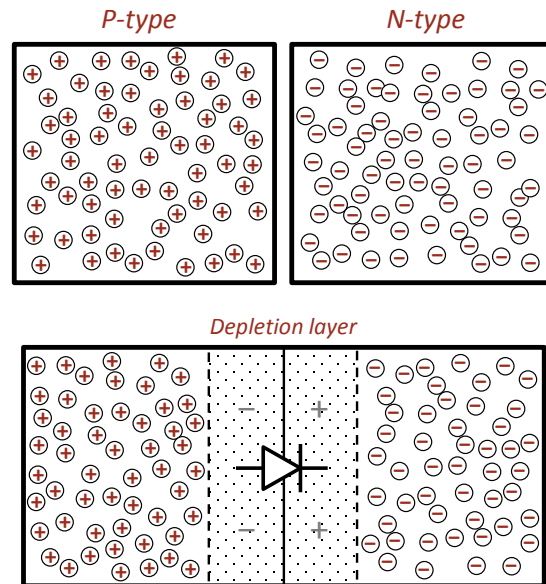


Figure 3-29: Formation of a PN junction diode and its depletion layer.

The PN junction I-V characteristic curve

The equilibrium configuration (Figure 3-29) is maintained when the rate that the holes diffuse into the depletion layer from the P-type side (against the contact potential gradient) matches the small rate of hole diffusion from the N-type side (where holes are the minority carriers), so that the net flow of holes across the junction is 0; a similar condition holds for the conduction electron diffusion at the junction.

Holes entering the depletion layer from the N-type side are not impeded by the presence of the contact potential — on the contrary, the electric field in the depletion layer will accelerate them through it to the P-type side. This implies that the rate of the minority hole diffusion will simply be proportional to the hole density on the N-type side, which is given in equation 3.9 to be proportional to $\exp(-q_e V_g/k_B T)$, and similarly for the minority electron diffusion from the P-type side. The majority carriers must cross the barrier imposed by the junction potential (V_j), so only those carriers with kinetic energies larger than $q_e V_j$ can cross to the other side; the number of such energetic carriers will be proportional to the Boltzmann factor $\exp(-q_e V_j/k_B T)$, because their kinetic energy distributions are classical, as mentioned before. At equilibrium, these two rates match, and $V_{j0} \approx V_g$.

When an external bias voltage V is applied across the PN junction, this applied potential will reduce the junction potential to $V_j = V_{j0} - V$ ($V > 0$ is forward-biased). As a consequence, more majority charge carriers will have enough energy to diffuse through the depletion layer; the minority diffusion rate from the other side is unaffected. Thus there will be a net current flow across the junction given by the difference in these two diffusion rates:

$$I \propto e^{-q_e(V_g - V)/k_B T} - e^{-q_e V_g/k_B T} = e^{-q_e V_g/k_B T} (e^{q_e V/k_B T} - 1)$$

This simple result is known as the *ideal diode equation*:

$$I = I_R (e^{q_e V/k_B T} - 1); \quad I_R = I_0 e^{-q_e V_g/k_B T}$$

V is the applied bias voltage (+ for forward-bias), V_g is the semiconductor's gap voltage, I_0 is some constant, and I_R is the diode's reverse leakage current. Thus, the ideal diode's forward current rises exponentially with forward bias voltage (for voltages of more than a few tens of millivolts), and has some small, temperature-dependent leakage current when reverse-biased.

The above equation is not quite right, because its derivation ignores an effect which is especially important for the behavior of a silicon diode: generation and recombination of charge carrier pairs in the depletion layer. The assumption in the argument leading up to the diode equation was that the only charge carriers present in the depletion layer entered it through diffusion from the regions outside the layer, and that all of these carriers pass through the depletion layer. Actually, thermal excitation of electron-hole pairs will occur in the depletion layer, just as it would in a pure semiconductor; similarly, recombination of electrons and holes may also occur among those which diffuse into the depletion layer, so the

number of charge carriers entering the depletion layer is larger than the number which escape, especially for small forward bias voltages in relatively large V_g diodes such as silicon.

The depletion layer generation and recombination processes depend exponentially on temperature, but the exponent goes as $q_e/2k_B T$ rather than as $q_e/k_B T$. The combination of this process with the ideal diode process leads to a “slight” modification of the ideal diode equation:

Diode equation

3.10

$$I = I_R(e^{q_e V/\eta k_B T} - 1); \quad I_R = I_0 e^{-q_e V_g/\eta k_B T}$$

The coefficient η depends on the importance of the depletion layer recombination process; it is a weak function of I and T and ranges between 1 and 2. For the small-signal silicon diodes you will use $\eta \approx 1.9$ and $I_R \approx 5 \text{ nA}$; $q_e/\eta k_B T \approx 20 \text{ volt}^{-1}$ at 20°C . The exponential dependence of I on forward-bias voltage V (for $I \gg I_R$) is the basis for the exponential and logarithmic amplifiers presented earlier. Unfortunately, according to (3.10) this relationship is strongly temperature-dependent. When a PN junction is reverse-biased, the thermal electron-hole generation rate in the depletion layer depends on the volume of the depletion layer, which grows as $\sqrt{1+V/V_g}$ (V is the reverse-bias voltage). This means that the reverse leakage current does not “saturate” at the I_R value given by equation (3.10), but continues to grow

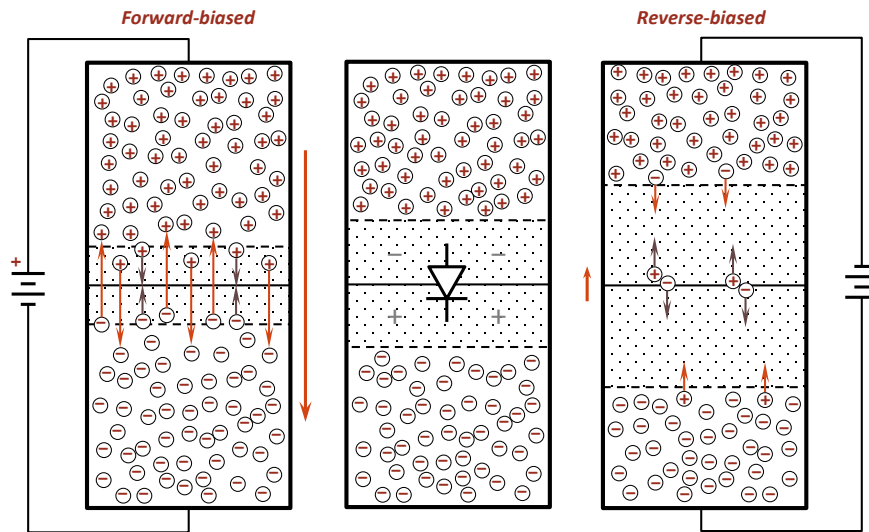


Figure 3-30: Depletion layer width and charge carrier diffusion of a PN junction as affected by applied bias voltage. Forward-bias (left) reduces the height of the potential barrier to majority carrier diffusion and decreases the depletion layer width, so many majority carriers can diffuse into and through the depletion layer; minority carrier diffusion is largely unaffected. Carriers that cross the depletion layer and recombine with majority carriers on the opposite side are indicated by the orange arrows; those that recombine inside the depletion layer are shown with gray arrows. Reverse-bias current is completely dominated by minority carrier diffusion: those that enter the depletion layer from the bulk semiconductor (orange arrows) and those pairs that are thermally generated within the depletion layer (gray arrows).

slowly with increasing reverse-bias voltage (as long as it remains well below the diode's reverse breakdown voltage, described in the next section). Cartoons of these diffusion processes are shown in Figure 3-30; plots of the 1N4148 silicon small-signal diode I-V characteristic curves for both forward and reverse bias and at two temperatures are provided in Figure 3-31.

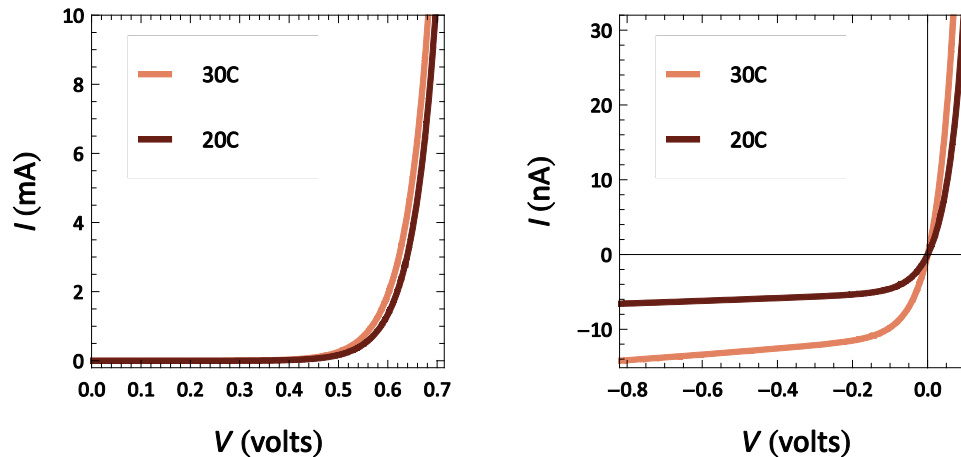


Figure 3-31: Forward-bias (left) and reverse-bias (right) I-V characteristic curves for the 1N4148 silicon diode. The forward-bias curves are exponential (equation 3.10), but they appear as though the diode suddenly “turns on” at a forward voltage of 0.6–0.7V; this voltage decreases slowly with rising diode temperature, as shown. The reverse-bias current is a much more sensitive function of temperature, however, approximately doubling for the same 10°C temperature increase. Note the different vertical scales for the two plots (a factor of 10^6).

Zener and avalanche breakdown

As the reverse-bias voltage on a PN junction is increased, the intensity of the electric field in the depletion layer rises; it is particularly intense at the interface between the P- and N-type areas. Minority carriers entering the depletion layer are accelerated by the field; when their kinetic energies reach a few eV or so, collisions with atoms in the lattice may knock valence electrons out of them, creating additional electron-hole pairs. These newly-created charge carriers are also accelerated by the field and can create even more carriers as they collide with lattice atoms.

At sufficiently high reverse voltages this collision-induced ionization process may lead to an *avalanche* of additional charge carriers, and the reverse current will grow exponentially with increasing voltage beyond some reverse-bias threshold. This is the *avalanche breakdown* process, and the reverse-bias voltage threshold for its action is the diode's *reverse breakdown voltage*. The electric field intensity for any particular applied reverse-bias voltage depends on the impurity concentrations and the abruptness with which these concentrations change near the P-type and N-type interface, so a target reverse breakdown voltage may be engineered into a particular diode type.

Another effect of a very intense electric field in the depletion layer is the large electric polarization of the atoms in the lattice it induces: at field strengths $\geq 10^6$ V/m the potential difference across a distance of about 100Å can exceed the semiconductor's gap voltage. In this case a valence electron may quantum mechanically *tunnel* across this distance into the conduction band, creating an electron-hole pair; because of this tunneling process the electric field required to ionize a lattice atom is much smaller than it would need to be to ionize a single, independent atom ($\sim 10^{10} - 10^{11}$ V/m); this effect was first theorized by the American physicist Clarence Zener in 1934. The tunneling rate grows exponentially as the required tunneling distance (inversely proportional to electric field strength) decreases, again leading to a large increase in reverse current (breakdown) as applied reverse-bias voltage exceeds the tunneling threshold. The target reverse breakdown voltage may be engineered by adjusting a diode's impurity concentration and doping profile. See Figure 3-20 on page 3-23 for a typical reverse breakdown I-V characteristic.

Diodes with reverse-breakdown voltages exceeding 6 V or so are dominated by the avalanche breakdown process; those below 5 V are predominantly subject to Zener breakdown. Regardless of breakdown voltage, those diodes designed to be used as voltage regulators with precisely-tailored reverse breakdown voltages are collectively called *Zener diodes*; those with high current-handling capacity and extremely fast response to voltages exceeding their breakdown threshold are usually called *avalanche diodes* and are primarily used for *transient voltage suppression* (TVS) and overvoltage protection.