



Chapter 1.2

Computational models of emotion

Stacy Marsella, Jonathan Gratch, and Paolo Petta

Summary

Recent years have seen a significant expansion in research on computational models of human emotional processes, driven both by their potential for basic research on emotion and cognition as well as their promise for an ever-increasing range of applications. This has led to a truly interdisciplinary, mutually beneficial partnership between emotion research in psychology and in computational science, of which this volume is an exemplar. To understand this partnership and its potential for transforming existing practices in emotion research across disciplines and for disclosing important novel areas of research, we explore in this chapter the history of work in computational models of emotion including the various uses to which they have been put, the theoretical traditions that have shaped their development, and how these uses and traditions are reflected in their underlying architectures.

For an outsider to the field, the last 15 years have seen the development of a seemingly bewildering array of competing and complementary computational models. Figure 1.2.1 lists a ‘family tree’ of a few of the significant models and the theoretical traditions from which they stem. Although there has been a proliferation of work, the field is far from mature: the goals that a model is designed to achieve are not always clearly articulated; research is rarely incremental, more often returning to motivating theories rather than extending prior computational approaches; and rarely are models contrasted with each other in terms of their ability to achieve their set goals. Contributing to potential confusion is the reality that computational models are complex systems embodying a number of, sometimes unarticulated, design decisions and assumptions inherited from the psychological and computational traditions from which they emerged, a circumstance made worse by the lack of a commonly accepted lexicon even for designating these distinctions.

In this chapter, we lay out the work on computational models of emotion in an attempt to reveal the common uses to which they may be put and the underlying techniques and assumptions from which the models are built. Our aim is to present conceptual distinctions and common terminology that can aid in discussion and comparison of competing models. Our hope is that this will not only facilitate an understanding of the field for outside researchers but work towards a lexicon that can help foster the maturity of the field towards more incremental research.

In characterizing different computational models of emotion, we begin by describing interdisciplinary uses to which computational models may be put, including their uses in improving human–computer interaction, in enhancing general models of intelligence, and as methodological tools for furthering our understanding of human behaviour. We next discuss how models have been built, including the underlying theoretical traditions that have shaped their development. These differing theoretical perspectives often conceptualize emotion in quite different ways, emphasizing different scenarios and proposed functions, different component processes, and different linkages between these components. It should then come as no surprise that such



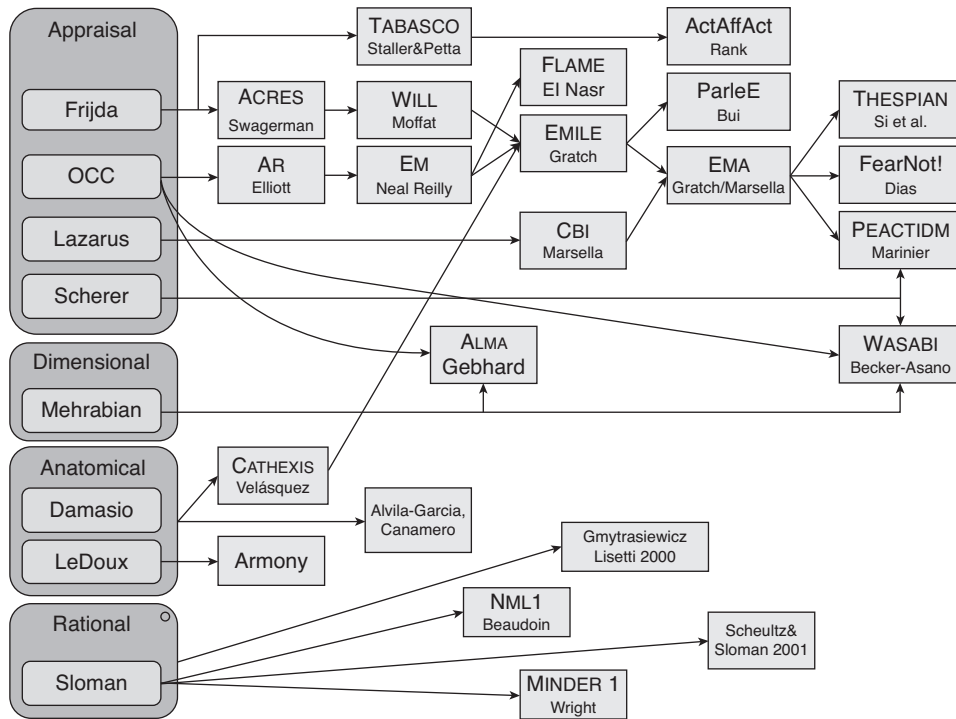


Fig. 1.2.1 A history of computational models of emotion.

1 differences are also reflected in the underlying design of the computational models. We next
 2 narrow our focus to cognitive appraisal theory, undeniably one of the most influential theoretical
 3 perspectives within computational research. To help organize and dissect research on computa-
 4 tional appraisal models, we introduce a generic appraisal architecture, a *component model view of*
 5 *appraisal models*, that conceptualizes emotion as a set of component models and relations between
 6 these components. We discuss how different computational systems address some, but typically
 7 not all, of these component models and describe differing processing choices that system develop-
 8 ers have used in realizing their component model variants. Finally, we illustrate how this compo-
 9 nent model view can help guide work in evaluating and contrasting alternative computational
 10 models of emotion.

11 The uses of computational models: an interdisciplinary partnership

12 New tools often transform science, opening up new approaches to research, allowing previously
 13 unaddressed questions to be explored, as well as revealing new questions. To appreciate the trans-
 14 formative role that computational models of emotion can have on research, we consider three
 15 aspects in this section: the impact on emotion research in psychology; the impact on artificial
 16 intelligence (AI); and, finally, the impact on work in human–computer interaction.

17 Impact on psychological research on emotion

18 Work in computational models of emotion impacts research in human emotion by transforming
 19 how theories are formulated and evaluated. One way this occurs is through a process of concretiz-
 20 ing concepts in the theory. Psychological theories of emotion have typically been cast at an
 21 abstract level and through informal (natural language) descriptions. Concepts in the theory are

usually not defined formally, and how processes work may not be laid out in systematic detail. The formulation of a computational model enforces more detail. The structures and processes of the theory must be explicitly and formally defined in order to implement them in a computational model, thus making a computer model a particularly concrete realization of the theory. The process of realizing the model can reveal implicit assumptions and hidden complexities, thereby forcing them to be addressed explicitly in some documented fashion. For example, appraisal theories often argue that a key variable in appraisal is an attribution of blameworthiness for an event deemed motivationally incongruent (e.g. Lazarus 1991). But the process by which a person makes such an attribution and therefore whether a particular situation would be deemed blameworthy, and the related required resources and capacities are typically not carefully laid out. And yet this attribution process may in itself be quite involved (e.g. Shaver 1985; Weiner 1995).

As computational modelling exposes hidden assumptions in the theory, addressing those assumptions can extend the scope of the theorizing. Seen in this way, computational models become not only a way to concretize theories, but also a framework for theory construction. In so doing, computational modelling also extends the language of emotion theorizing by incorporating concepts, processes, and metaphors drawn from computation, much as concepts such as *information processing* and *symbol systems* have impacted psychology in general. For example, several computational models have recast the appraisal theory in terms of concepts drawn from AI, including knowledge representation (e.g. Gratch and Marsella 2004a), planning (e.g. Gratch 2000; Dias and Paiva 2005), neural networks (Sander *et al.* 2005), and decision making (Lisetti and Gmytrasiewicz 2002; Ito *et al.* 2008). Incorporation of the models into larger simulations can also expose hidden questions behind traditional conceptualizations and extend the scope of theorizing. For example, several computational models of emotion have been incorporated into larger simulation systems that seek to model emotion's role in human mental processes and behaviour (Marsella and Gratch 2001; Dias and Paiva 2005; Becker-Asano 2008; Rank 2009). This has led researchers to address fundamental architectural questions about the relation of appraisal processes to other cognitive processes, perception, and behaviour (Marsella and Gratch 2009; Rank 2009). Of course, a central challenge here is to ensure that increases in the scope of the theorizing do not endanger the parsimony often critical to a model's explanatory power.

Coupled to this transformation of the theory formation process through modelling and simulation runs of the model, the computational realization of a theory can also increase the capacity to draw predictions from theory. In particular, computational models provide a new empirical framework for studying emotion processes that goes beyond what is feasible in more traditional laboratory settings. Computer simulations of the model *behave*: they provide a means to explore systematically the temporal dynamics of emotion processes and form predictions about the time course of those processes. Manipulations of experimental conditions may be explored more extensively first with a computational model, such as ablating certain functionalities or testing responses under adverse conditions that may be costly, risky, or raise ethical concerns *in vivo* (e.g. Armony *et al.* 1997). Simulations can reveal unexpected model properties that suggest further exploration. Additionally, models of emotion and affective expression have been incorporated into virtual humans, software artefacts that look and act like humans, capable of perception and action in a virtual world that they can cohabit with people. These systems essentially allow for the study of emotion in a virtual ecology, a form of synthetic *in vivo* experimentation.

Finally, the computational modelling of emotion and emotional expression has led to new ways to create stimuli for human subject experimentation. Virtual humans are in some ways the experimenter's ultimate confederate. A virtual human can be manipulated systematically to elicit behaviour from human subjects. For example, virtual humans have been used to show that subtle

1 changes in physical appearance or behaviour can profoundly impact social interaction, including
 2 changes to people's willingness to cooperate (Krumhuber *et al.* 2007), the fluidity of their conver-
 3 sation (Gratch *et al.* 2007), learning outcomes (Baylor and Kim 2008) and even their level of
 4 social aggression (McCall *et al.* 2009).

5 **Impact on artificial intelligence and robotics**

6 Modern research in the psychology, cognitive science, and neuroscience of emotion has led to a
 7 revolution in our thinking about the relation of emotion to cognition and social behaviour and,
 8 as a consequence, is also transforming the science of computation. Findings on the functional,
 9 often adaptive role that emotions play in human behaviour have motivated AI and robotics
 10 research to explore whether computer analogues of human emotion can lead to more intelligent,
 11 flexible, and capable systems. Early work by Simon (1967) argued that emotions serve the crucial
 12 function of interrupting normal cognition when unattended goals require servicing. Viewing
 13 emotion as serving this critical interrupt capacity provides a means for an organism to balance
 14 competing goals as well as incorporate reactive behaviours into more deliberative processing.
 15 A range of studies point to emotions as the means by which the individual establishes values for
 16 alternative decisions and decision outcomes. Busemeyer *et al.* (2007) argue that emotional state
 17 influences the subjective utility of alternative choice. Studies performed by Damásio and col-
 18 leagues suggest that damage to the ventromedial prefrontal cortex prevents emotional signals
 19 from guiding decision making in an advantageous direction (Bechara *et al.* 1999).

20 Other authors have emphasized how social emotions such as anger and guilt may reflect a
 21 mechanism that improves group utility by minimizing social conflicts, and thereby explains peo-
 22 ple's 'irrational' choices to cooperate in social games such as prisoner's dilemma (Frank 1988).
 23 Similarly, 'emotional biases' such as wishful thinking may reflect a rational mechanism that more
 24 accurately accounts for certain social costs, such as the cost of betrayal when a parent defends a
 25 child despite strong evidence of their guilt in a crime (Mele 2001).

26 Collectively, these findings suggest that emotional influences have important social and
 27 cognitive functions that would be required by *any* intelligent system. This view is not new to arti-
 28 ficial intelligence (Simon 1967; Sloman and Croucher 1981; Minsky 1986) but was in large meas-
 29 ure ignored in AI research of the late twentieth century which largely treated emotion as
 30 antithetical to rationality and intelligence. However, in the spirit of Hume's famous dictum: 'rea-
 31 son is, and ought only to be the slave of the passions...' (Hume 2000, 2.3.3.4), the question of
 32 emotion has again come to the fore in AI as models have begun to catch up with theoretical find-
 33 ings. This has been spurred, in part, by an explosion of interest in integrated computational
 34 models that incorporate a variety of cognitive functions (Bates *et al.* 1991; Anderson 1993; Rickel
 35 *et al.* 2002). Indeed, until the rise of broad integrative models of cognition, the problems emotion
 36 was purported to solve, for example, juggling multiple goals, were largely hypothetical. More
 37 recent cognitive systems embody a variety of mental functions and face very real challenges how
 38 to allocate resources. A recurring theme in emotion research in AI is the role of emotion in
 39 addressing such control choices by directing cognitive resources towards problems of adaptive
 40 significance for the organism (Scheutz and Sloman 2001; Staller and Petta 2001; Blanchard and
 41 Cañamero 2006; Scheutz and Schermerhorn 2009).

42 **Impact on human-computer interaction**

43 Finally, research has revealed the powerful role that emotion and emotion expression play in
 44 shaping human social interaction, and this in turn has suggested that computer interaction
 45 can exploit (and indeed must address) this function. Emotional displays convey considerable

1 information about the mental state of an individual. Although there is a lively debate about
 2 whether these displays reflect true emotion or are simply communicative conventions (Manstead
 3 *et al.* 1999), pragmatically there is truth in both perspectives. From emotional displays, observers
 4 can form interpretations of a person's beliefs (e.g. frowning at an assertion may indicate disagree-
 5 ment), desires (e.g. joy gives information that a person values an outcome), and intentions/action
 6 tendencies (e.g. fear suggests flight). They may also provide information about the underlying
 7 dimensions along which people appraise the emotional significance of events: valence; intensity;
 8 certainty; expectedness; blameworthiness; etc. (Smith and Scott 1997). With such a powerful
 9 signal, it is not surprising that emotions can be a means of social control (Fridlund 1997; Campos
 10 *et al.* 2003; de Waal 2003). Emotional displays seem to function to elicit particular social responses
 11 from other individuals ('social imperatives', Frijda 1987) and arguably, such responses can be
 12 difficult to suppress. The responding individual may not even be consciously aware of the manip-
 13 ulation. For example, anger seems to be a mechanism for coercing actions in others and enforcing
 14 social norms, displays of guilt can elicit reconciliation after some transgression, distress can be
 15 seen as a way of recruiting social support, and displays of joy or pity are a way of signalling such
 16 support to others. Other emotion displays seem to exert control indirectly, by inducing emotional
 17 states in others and thereby influencing an observer's behaviour. Specific examples of this are
 18 *emotional contagion*, which can lead individuals to 'catch' the emotions of those around them
 19 (Hatfield *et al.* 1994) and the *Pygmalion effect* (also known as 'self-fulfilling prophecy') whereby
 20 our positive or negative expectations about an individual, even if expressed nonverbally, can
 21 influence them to meet these expectations (Blanch 1993). Given this wide array of functions in
 22 social interactions, many have argued that emotions evolved because they provide an adaptive
 23 advantage to social organisms (Darwin 1872/1998; de Waal 2003).

24 To the extent that these functions can be realized in artificial systems, they could play a
 25 powerful role in facilitating interactions between computer systems and human users. This has
 26 inspired several trends in human-computer interaction. For example, Conati uses a Bayesian
 27 network-based appraisal model to deduce a student's emotional state based on their actions
 28 (Conati 2002) and several systems have attempted to recognize the behavioural manifestations of
 29 a user's emotion including facial expressions (Lisetti and Schiano 2000; Fasel *et al.* 2002; Haag
 30 *et al.* 2004), physiological indicators (Picard 1997; Haag *et al.* 2004), and vocal expression (Lee
 31 and Narayanan 2003).

32 A related trend in human-computer interaction (HCI) work is the use of emotions and
 33 emotional displays in virtual characters that interact with the user. As animated films (Thomas
 34 and Johnston 1995) so poignantly demonstrate, emotional displays in an artificially generated
 35 character can have the general effect of making it seem human or lifelike, and thereby cue the user
 36 to respond to, and interact with, the character as if it were another person. A growing body of
 37 research substantiates this view. In the presence of a lifelike agent, people are more polite, tend to
 38 make socially desirable choices, and are more nervous (Kramer *et al.* 2003); they can exhibit
 39 greater trust of the agent's recommendations (Cowell and Stanney 2003) and they can feel more
 40 empathy (Paiva *et al.* 2004). In that people utilize these behaviours in their everyday interpersonal
 41 interactions, modelling the function of these behaviours is essential for any application that hopes
 42 to faithfully mimic face-to-face human interaction. More importantly, however, the ability of
 43 emotional behaviours to influence a person's emotional and motivational state could potentially,
 44 if exploited effectively, guide a user towards more effective interactions. For example, education
 45 researchers have argued that nonverbal displays can have a significant impact on student intrinsic
 46 motivation (Lepper 1988).

47 A number of applications have attempted to exploit this interpersonal function of emotional
 48 expression. Klesen (2005) models the communicative function of emotion, using stylized

1 animations of body language and facial expression to convey a character's emotions and inten-
 2 tions with the goal of helping students understand and reflect on the role these constructs play in
 3 improvisational theater. Nakanishi *et al.* (2005) and Cowell and Stanney (2003) each evaluated
 4 how certain nonverbal behaviours could communicate a character's trustworthiness for training
 5 and marketing applications, respectively. Several applications have also tried to manipulate a
 6 student's motivations through the emotional behaviours of a virtual character. Lester utilized
 7 praising and sympathetic emotional displays to provide feedback and increase student motivation
 8 in a tutoring application (Lester *et al.* 2000). Researchers have also looked at emotion and emo-
 9 tional expression in characters as a means to engender empathy and bonding between between
 10 learners and virtual characters (Marsella *et al.* 2003; Paiva *et al.* 2005); Biswas and colleagues
 11 (Biswas *et al.* 2005) also use human-like traits to promote empathy and intrinsic motivation in a
 12 learning-by-teaching system.

13 In summary, computational models of emotion serve differing roles in research and applica-
 14 tions. Further, the evaluation of these models is in large measure dependent on those roles. In the
 15 case of the psychological research that uses computational models, the emphasis will largely be on
 16 fidelity with respect to human emotion processes. In the case of work in AI and robotics, evalua-
 17 tion often emphasizes how the modelling of emotion impacts reasoning processes or leads in
 18 some way to improved performances such as an agent or robot that achieves a better fit with its
 19 environment. In HCI work, the key evaluation is whether the model improves human-computer
 20 interaction such as by making it more effective, efficient, or pleasant.

21 Overall, the various roles for computational models of emotion have led to a number of
 22 impressive models being proposed and developed. To put this body of work into perspective, it
 23 is critical for the field to support a deeper understanding of the relationship between these mod-
 24 els. To assist in that endeavour, we now turn to presenting some common terms and distinctions
 25 that can aid in the comparison of competing models.

26 **A component perspective on the design of computational models**

27 Each of the computational models listed in Figure 1.2.1 is a very different entity, with incompat-
 28 ible inputs and outputs, different behaviours, embodying irreconcilable processing assumptions
 29 and directed towards quite different scientific objectives. What we argue here, however, is that
 30 much of this variability is illusory. These models are complex systems that integrate a number of
 31 component 'submodels'. Sometimes these components are not clearly delineated, but, if one dis-
 32 assembles models along the proper joints, then a great many apparent differences collapse into a
 33 small number of design choices. To facilitate this decomposition, this section describes the com-
 34 ponent processes underlying emotion, with a particular emphasis on components posited in con-
 35 nection with appraisal theory. These components are not new—indeed, they are central
 36 theoretical constructs in many theories of emotion—but some of the terminology is new as we
 37 strive to simplify terms and de-conflict them with other terminology more commonly used in
 38 computer science. We begin by describing the various theoretical traditions that have influenced
 39 computational research and the components these theories propose.

40 A challenge in developing a coherent framework for describing computational models of
 41 emotion is that the term 'emotion' itself is fraught with ambiguities and contrasting definitions.
 42 Emotions are a central aspect of everyday life and people have strong intuitions about them. As a
 43 consequence, the terms used in emotion research (appraisal, emotion, mood, affect, feeling) have
 44 commonsense interpretations that can differ considerably from their technical definition within
 45 the context of a particular emotion theory or computational model (Russell 2003). This ambigu-
 46 ity is confounded by the fact that there are fundamental disputes within psychological and

1 neuroscience research on emotion over the meaning and centrality of these core concepts.
 2 Theories differ as to which components are intrinsic to an emotion (e.g. cognitions, somatic
 3 processes, behavioural tendencies and responses), the relationships between components (e.g. do
 4 cognitions precede or follow somatic processes), and representational distinctions (e.g. is anger a
 5 linguistic fiction or a natural kind)—see Chapter 1.1, this volume, for an overview of different
 6 theoretical perspectives on emotion.

7 Understanding these alternative theoretical perspectives on emotion is essential for anyone
 8 who aspires to develop computational models, but this does not imply that a modeller must be
 9 strictly bound by any specific theoretical tradition. Certainly, modellers should strive for a con-
 10 sistent and well-founded semantics for their underlying emotional constructs and picking and
 11 integrating fundamentally irreconcilable theoretical perspectives into a single system can be prob-
 12 lematic at best. If the goal of the computational model is to faithfully model human emotional
 13 processes or, more ambitiously, to contribute to theoretical discourse on emotion, such incon-
 14 sistencies can be fatal. However, some ‘fundamentally irreconcilable’ differences are illusory and
 15 evaporate when seen from a new perspective. For example, disputes as to whether emotion pre-
 16 cededs or follows cognition dissipate if one adopts a dynamic systems perspective (i.e. a circle has
 17 no beginning). Nonetheless, theoretical models provide important insights in deriving a coherent
 18 computational model of emotion and deviations from specific theoretical constraints, ideally,
 19 will be motivated by concrete challenges in realizing a theory within a specific representational
 20 system or in applying the resulting model to concrete applications. Here we review some of the
 21 theoretical perspectives that have most influenced computational modelling research.

22 Appraisal theory

23 Appraisal theory, discussed in detail in Chapter 1.1, is currently a predominant force among
 24 psychological perspectives on emotion and arguably the most fruitful source for those interested
 25 in the design of symbolic AI systems, as it emphasizes and explains the connection between emo-
 26 tion and cognition. Indeed, the large majority of computational models of emotion stem from
 27 this tradition. In appraisal theory, emotion is argued to arise from patterns of individual judge-
 28 ment concerning the relationship between events and an individual’s beliefs, desires, and inten-
 29 tions, sometimes referred to as the *person–environment relationship* (Lazarus 1991). These
 30 judgements, formalized through reference to devices such as *situational meaning structures*
 31 or *appraisal variables* (Frijda 1987), characterize aspects of the personal significance of events.
 32 Patterns of appraisal are associated with specific physiological and behavioural reactions. In sev-
 33 eral versions of appraisal theory, appraisals also trigger cognitive responses, often referred to as
 34 *coping strategies*—e.g. planning, procrastination, or resignation—feeding back into a continual
 35 cycle of appraisal and reappraisal (Lazarus 1991, p. 127).

36 In terms of underlying components of emotion, appraisal theory foregrounds appraisal as a
 37 central process. Appraisal theorists typically view appraisal as the cause of emotion, or at least of
 38 the physiological, behavioural, and cognitive changes associated with emotion. Some appraisal
 39 theorists emphasize ‘emotion’ as a discrete component within their theories, whereas others treat
 40 the term emotion more broadly to refer to some configuration of appraisals, bodily responses,
 41 and subjective experience (see Ellsworth and Scherer 2003 for a discussion). Much of the work
 42 has focused on the structural relationship between appraisal variables and specific emotion
 43 labels—e.g. which pattern of appraisal variables would elicit hope (see Ortony *et al.* 1988)—or
 44 the structural relationship between appraisal variables and specific behavioural and cognitive
 45 responses—e.g. which pattern of appraisal variables would elicit certain facial expressions
 46 (Smith and Scott 1997; Scherer and Ellgring 2007a) or coping tendencies (Lazarus 1991). Indeed,
 47 although appraisal theorists allow that the same situation may elicit multiple appraisals, theorists

are relatively silent on how these individual appraisals would combine into an overall emotional state or if this state is best represented by discrete motor programs or more dimensional representations. More recent work has begun to examine the processing constraints underlying appraisal—to what extent is it parallel or sequential (Scherer 2001; Moors *et al.* 2005) and does it occur at multiple levels (Smith and Kirby 2000; Scherer 2001)—and creating a better understanding of the cognitive, situational, and dispositional factors that influence appraisal judgements (Kuppens and Van Mechelen 2007; Smith and Kirby 2009).

Models derived from appraisal theories of emotion, not surprisingly, emphasize appraisal as the central process to be modelled. Computational appraisal models often encode elaborate mechanisms for deriving appraisal variables such as decision-theoretic plans (Gratch and Marsella 2004a; Marsella and Gratch 2009), reactive plans (Staller and Petta 2001; Rank and Petta 2005; Neal Reilly 2006), Markov-decision processes (El Nasr *et al.* 2000; Si *et al.* 2008), or detailed cognitive models (Marinier *et al.* 2009). Emotion itself is often less elaborately modelled. It is sometimes treated simply as a label (sometimes with an intensity) to which behaviour can be attached (Elliott 1992). Appraisal is typically modelled as the cause of emotion with the specific emotion label being derived via *if-then rules* on a set of appraisal variables. Some approaches make a distinction between a specific emotion instance (allowing multiple instances to be derived from the same event) and a more generalized ‘affective state’ or ‘mood’ (see discussion of *core affect* below) that summarizes the effect of recent emotion elicitations (Neal Reilly 1996; Gratch and Marsella 2004a; Gebhard 2005). Some more recent models attempt to model the impact of momentary emotion and mood on the appraisal process (Gratch and Marsella 2004a; Gebhard 2005; Paiva *et al.* 2005; Marsella and Gratch 2009).

Computational appraisal models have been applied to a variety of uses including contributions to psychology, AI, and HCI. For example, Marsella and Gratch have used EMA to generate specific predictions about how human subjects will appraise and cope with emotional situations and argue that empirical tests of these predictions have implications for psychological appraisal theory (Gratch, *et al.* 2009b; Marsella *et al.* 2009). Several authors have argued that appraisal processes would be required by any intelligent agent that must operate in real-time, ill-structured, multi-agent environments (e.g. Staller and Petta 2001). The bulk of application of these techniques, however, has been for HCI applications, primarily for the creation of real-time interactive characters that exhibit emotions in order to make these characters more compelling (e.g. Neal Reilly 1996), more realistic (e.g. Traum *et al.* 2003; Mao and Gratch 2006), more able to understand human motivational state (e.g. Conati and MacLaren 2004) or more able to induce desirable social effects in human users (e.g. Paiva *et al.* 2005).

Dimensional theories

Dimensional theories of emotion argue that emotion and other affective phenomena should be conceptualized, not as discrete entities but as points in a continuous (typically two- or three-) dimensional space (Mehrabian and Russell 1974; Watson and Tellegen 1985; Russell 2003; Barrett 2006). Indeed, many dimensional theories argue that discrete emotion categories (e.g. hope, fear, and anger) are folk-psychological concepts that have unduly influenced scientific discourse on emotion and have no ‘reality’ in that there are no specific brain regions or circuits that correspond to specific emotion categories (Barrett 2006). Not surprisingly, dimensional theories de-emphasize the term emotion or relegate it to a cognitive label attributed, retrospectively, to some perceived body state. Rather they emphasize concepts such as mood, affect, or, more recently, *core affect* (Russell 2003). We adopt this later term in subsequent discussion. A person is said to be in exactly one affective state at any moment (Russell 2003, p. 154) and the space of

possible core affective states is characterized in terms of broad, continuous dimensions. Many computational dimensional models build on the three-dimensional 'PAD' model of Mehrabian and Russell (1974) where these dimensions correspond to *pleasure* (a measure of valence), *arousal* (indicating the level of affective activation), and *dominance* (a measure of power or control).

It is worth noting that there is a relationship between the dimensions of core affect and appraisal dimensions—the pleasure dimension roughly maps on to appraisal dimensions that characterize the valence of an appraisal-eliciting event (e.g. intrinsic pleasantness or goal congruence), dominance roughly maps on to the appraisal dimension of coping potential, and arousal is a measure of intensity. However, they have quite different meanings: appraisal is a relational construct characterizing the relationship between some specific object/event and the individual's beliefs, desires, and intentions and several appraisals may be simultaneously active; core affect is a non-relational construct summarizing a unique overall state of the individual.

Dimensional theories emphasize different components of emotion than appraisal theories and link these components quite differently. Dimensional theories foreground the structural and temporal dynamics of core affect and often do not address affect's antecedents in detail. Most significantly, dimensional theorists question the tight causal linkage between appraisal and emotion that is central to appraisal accounts. Dimensional theorists conceive of core affect as a 'non-intentional' state, meaning the affect is not about some object (as in 'I am angry at *him*'). In such theories, many factors may contribute to a change in core affect including symbolic intentional judgements (e.g. appraisal) but also subsymbolic factors such as hormones and drugs (Schachter and Singer 1962), but, most importantly, the link between any preceding intentional meaning and emotion is broken (as it is not represented within core affect) and must be recovered after the fact, sometimes incorrectly (Clore *et al.* 1994; Clore and Plamer 2009). For example, Russell argues for the following sequence of emotional components: some external event occurs (e.g. a bear walks out of the forest), it is perceived in terms of its affective quality; this perception results in a dramatic change in core affect; this change is attributed to some 'object' (e.g. the bear); and only then is the object cognitively appraised in terms of its goal relevance, causal antecedents, and future prospects (see also Zajonc 1980).

Models influenced by dimensional theories, not surprisingly, emphasize processes associated with core affect and other components (e.g. appraisal) tend to be less elaborately developed. Core affect is typically represented as a continuous time-varying process that is represented at a given period of time by a point in three-dimensional space that is 'pushed around' by eliciting events. Computational dimensional models often have detailed mechanisms for how this point changes over time—e.g. decays to some resting state—and incorporating the impact of dispositional tendencies such as personality or temperament (Gebhard 2005).

Computational dimensional models are most often used for animated character behaviour generation, perhaps because it translates emotion into a small number of continuous dimensions that can be readily mapped to continuous features of behaviour such as the spatial extent of a gesture. For example, PAD models describe all behaviour in terms of only three dimensions, whereas modellers using appraisal models must either associate behaviours with a larger number of appraisal dimensions (see Smith and Scott 1997; Scherer and Ellgring 2007a) or map appraisals into a small number of discrete, though perhaps intensity-varying, expressions (Elliott 1992). For a similar reason, dimensional models are also frequently used as a good representational framework for systems that attempt to recognize human emotional behaviour and there is some evidence that they may better discriminate user affective states than approaches that rely on discrete labels (Barrett 2006).

The relationship between core affect and cognition is generally less explored in dimensional approaches. Typically the connection between emotion-eliciting events and current core-affective

1 state is not maintained, consistent with Russell's view of emotion as a non-intentional state
 2 (e.g. Becker-Asano and Wachsmuth 2008; see Chapter 4.1, this volume). Interestingly, we are not
 3 aware of any computational models that follow the suggestion from Zajonc and Russell that
 4 appraisal is a *post hoc* explanation of core affect. Rather, many computational models of emotion
 5 that incorporate core affect have viewed appraisal as the mechanism that initiates changes to core
 6 affect. For example, Gebhard's (2005) ALMA model includes Ortony, *et al.* (1988) inspired
 7 appraisal rules (OCC) and WASABI (Becker-Asano and Wachsmuth 2008; see Chapter 4.1, this
 8 volume) incorporates appraisal processes inspired by Scherer's sequential-checking theory into a
 9 PAD-based model of core affect. Some computational models explore how core affect can influ-
 10 ence cognitive processes. For example, HOTCO 2 (Thagard 2003) allows explanations to be
 11 biased by dimensional affect (in this case, a one-dimensional model encoding valence) but this
 12 is more naturally seen as the consequence of emotion on cognition (e.g. the modelling of
 13 an emotion-focused coping strategy in the sense of Lazarus 1991).

14 Other approaches

- 15 ♦ *Anatomic theories* stem from an attempt to reconstruct the neural links and processes that
 16 underlie organisms' emotional reactions (Le Doux 1996; Panksepp 1998a; Öhman and Wiens
 17 2004). Unlike appraisal theories, such models tend to emphasize subsymbolic processes.
 18 Unlike dimensional theories, anatomic approaches tend to view emotions as different, dis-
 19 crete neural circuits and emphasize processes or systems associated with these circuits. Thus,
 20 anatomically inspired models tend to foreground certain process assumptions and tend to be
 21 less comprehensive than either appraisal or dimensional theories, with researchers focusing
 22 on a specific emotion such as fear. For example, LeDoux emphasizes a 'high-road' versus 'low-
 23 road' distinction in the fear circuit with the latter reflecting automatic/reflexive responses
 24 to situations, whereas the former is mediated by cognition and deliberation. Computational
 25 models inspired by the anatomic tradition often focus on low-level perceptual-motor
 26 tasks and encode a two-process view of emotion that argues for a fast, automatic, undifferenti-
 27 ated emotional response and a slower, more differentiated response that relies on higher-level
 28 reasoning processes (e.g. Armony *et al.* 1997).
- 29 ♦ *Rational approaches* start from the question of what adaptive function does emotion serve and
 30 then attempt to abstract this function away from its 'implementation details' in humans and
 31 incorporate these functions into a (typically normative) model of intelligence (Simon 1967;
 32 Sloman and Croucher 1981; Frank 1988; Scheutz and Sloman 2001; Anderson and Lebiere
 33 2003; Doyle 2006). Researchers in this tradition typically reside in the field of artificial intelli-
 34 gence and view emotion as a window through which one can gain insight into adaptive behav-
 35 iour, albeit it a very different window than that which has motivated much of AI research.
 36 Within this tradition, cognition is conceived as a collection of symbolic processes that serve
 37 specific cognitive functions and are subject to certain architectural constraints on how they
 38 interoperate. Emotion, within this view, is simply another, albeit often overlooked, set of
 39 processes and constraints that have adaptive value. Models of this sort are most naturally
 40 directed towards the goal of improving theories of machine intelligence.
- 41 ♦ *Communicative theories of emotion* argue that emotion processes function as a communicative
 42 system: both as a mechanism for informing other individuals of one's mental state—and
 43 thereby facilitate social coordination—and as a mechanism for requesting/demanding changes
 44 in the behaviour of others—as in threat displays (Keltner and Haidt 1999; Parkinson 2009).
 45 Communicative theories emphasize the social-communicative function of displays and some-
 46 times argue for a disassociation between internal emotional processes and emotion displays

that need not be selected on the basis of an internal emotional state (e.g. see Fridlund 1997; Gratch 2008). Computational models inspired by communicative theories often embrace this disassociation and dispense with the need for an internal emotional model and focus on machinery that decides when an emotional display will have a desirable effect on a human user. For example, in the Cosmo tutoring system (Lester, Towns, *et al.* 2000), the agent's pedagogical goals drive the selection and sequencing of emotive behaviours. In Cosmo, a congratulatory act triggers a motivational goal to express admiration that is conveyed with applause. Not surprisingly, computational models based on communicative theories are most often directed towards the goal of achieving social influence.

Dissecting computational appraisal theory

Appraisal theory, by far, dominates the work on computational models of emotion so here we spend some time laying out some terminology that is specific to this class of models (although some of this terminology could apply to other approaches). As we discussed earlier, our aim is to promote incremental research on computational models of emotion by presenting a compositional view of model building, emphasizing that an emotional model is often assembled from individual 'submodels' and these smaller components are often shared and can be mixed, matched, or excluded from any given implementation. More importantly, these components can be seen as embodying certain content and process assumptions that can be potentially assessed and subsequently abandoned or improved as a result of these assessments. In presenting this, we attempt to build as much as possible on terminology already introduced within the emotion literature.

Figure 1.2.2 presents an idealized computational appraisal architecture consisting of a set of linked component models. This figure presents what we see as natural joints at which to decompose appraisal systems into coherent and often shared modules, though any given system may fail to implement some of these components or allow different information paths between components. In this architecture, information flows in a cycle as argued by several appraisal theorists (Lazarus 1991; Scherer 2001; Parkinson 2009): some representation of the person–environment relationship is appraised; this leads to an affective response of some intensity; the response triggers behavioural and cognitive consequences; these consequences alter the person–environment;

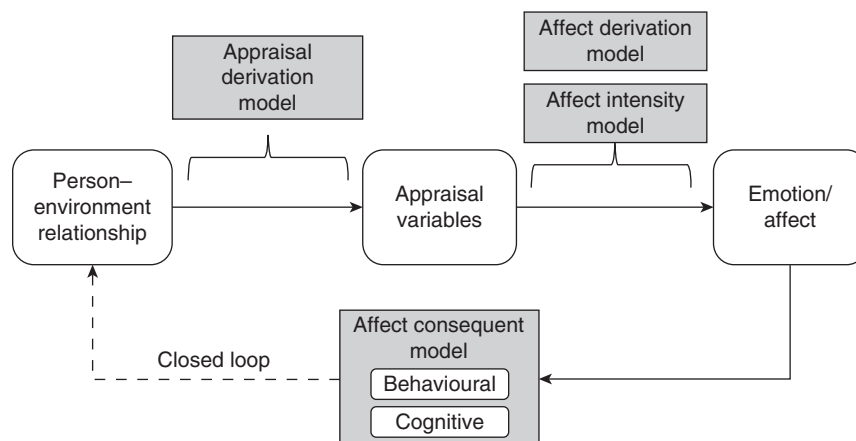


Fig. 1.2.2 A component model view of computational appraisal models.

1 this change is appraised; and so on. Each of these stages can be represented by a model that
 2 represents or transforms state information relevant to emotion-processing. Here we introduce
 3 terminology associated with each of these.

4 **Person–environment relationship**

5 Lazarus (1991) introduced this term to refer to some representation of the agent's relationship
 6 with its environment. This representation should allow an agent, in principle, to derive the rela-
 7 tionship between external events (real or hypothetical) and the beliefs, desires, and intentions
 8 of the agent or other significant entities in the (real or hypothetical) social environment.
 9 This representation need not represent these relationships explicitly but must support their deri-
 10 vation. Examples of this include the decision-theoretical planning representations in EMA
 11 (Gratch and Marsella 2004), which combine decision-theoretic planning representation with
 12 belief–desire–intention formalisms, or the partially observable Markov-decision representations
 13 in THESPIAN (Si *et al.* 2008)

14 **Appraisal-derivation model**

15 An appraisal-derivation model transforms some representation of the person–environment
 16 relationship into a set of appraisal variables.¹ For example, if an agent's goal is potentially thwarted
 17 by some external action, an appraisal-derivation model should be able to automatically derive
 18 appraisals that this circumstance is undesirable, assess its likelihood, and calculate the agent's
 19 ability to cope, i.e. by identifying alternative ways to achieve this goal. Several computational
 20 appraisal models don't provide an appraisal-derivation model or else treat its specification as
 21 something that is outside of the system. For example, ALMA (Gebhard 2005) allows domain
 22 developers to author the relational model by specifying how certain states or actions should be
 23 appraised (e.g. if Sven attacks Valerie, she should appraise this as undesirable). Other researchers
 24 treat the appraisal-derivation as a central contribution of their approach. For example, EMA
 25 provides a series of domain-independent inference rules that derive appraisal variables from syn-
 26 tactic features of the person–environment relationship (e.g. if the effect of an action threatens a
 27 plan to achieve a desired state, this is undesirable). Models also differ in the processing constraints
 28 that this model should satisfy. For example, models influenced by Scherer's sequential checking
 29 theory incorporate assumptions about the order in which specific appraisal variables should be
 30 derived (Marinier 2008). Appraisal-derivation models are often triggered by some eliciting event,
 31 though this is not always the case (e.g. EMA simultaneously appraises every goal in an agent's
 32 working memory and updates these appraisals continuously as new information about these
 33 goals is obtained.

34 **Appraisal variables**

35 Appraisal variables correspond to the set of specific judgements that the agent can use to produce
 36 different emotional responses and are generated as a result of an appraisal-derivation model.
 37 Different computational appraisal models adopt different sets of appraisal variables, depending
 38 on their favourite appraisal theorist. For example, many approaches utilize the set of variables
 39 proposed by Ortony, *et al.* (1988) including AR (Elliott 1992), EM (Neal Reilly 1996), FLAME

¹ Smith and Kirby (2009) propose the term *relational model* to refer to this mapping, building on Lazarus's idea that appraisal is a relational construct relating the person and the environment. They introduced the term to draw attention to the fact that many appraisal theories emphasize the mapping from appraisal variable to emotion but neglect the situational and dispositional antecedents of appraisal. As 'relation' and 'relational' often have a very different meaning within computer science, we prefer a different term.

1 (El Nasr, *et al.* 2000), and ALMA (Gebhard 2005). Others favour the variables proposed by Scherer
 2 (2001) including WASABI (Becker-Asano and Wachsmuth 2008; see Chapter 4.1, this volume)
 3 and PEACTION (Marinier *et al.* 2009).

4 Affect-derivation model

5 An affect-derivation model maps between appraisal variables and an affective state, and specifies
 6 how an individual will react emotionally once a pattern of appraisals has been determined.² As
 7 noted in the discussion of different theoretical perspectives above, there is some diversity in how
 8 models define ‘emotion’ and here we consider any mapping from appraisal variables to affective
 9 state, where this state could be either a discrete emotion label, a set of discrete emotions, core
 10 affect, or even some combination of these factors. For example, Elliott’s AR (Elliott 1992) maps
 11 appraisal variables into discrete emotion labels, Becker-Asano’s WASABI (Becker-Asano and
 12 Wachsmuth 2008) maps appraisals into a dimensional (e.g. PAD) representation of emotion,
 13 and Gebhard’s ALMA (Gebhard 2005) does both simultaneously. Many computational systems
 14 adopt the affect-derivation model proposed by Ortony, *et al.* (1988) whereby 22 emotion labels
 15 are defined as conjunctions of appraisal variables—this will be henceforth referred to as the OCC
 16 model. Others have implemented models based on the work of Lazarus (e.g. Gratch and Marsella’s
 17 EMA) and Scherer (e.g. Becker-Asano’s WASABI and Marinier’s PEACTION). Much of
 18 the empirical work in psychological appraisal theory has focused on identifying the affect-
 19 derivation model that best conforms to human behaviour but the results of these studies are far
 20 from definitive and can be interpreted to support multiple proposed models.

21 Affect-intensity model

22 An affect-intensity model specifies the strength of the emotional response resulting from a
 23 specific appraisal. There is a close association between the affect-derivation model and intensity
 24 model (the intensity computation is often implemented as part of the affect-derivation model).
 25 However it is useful to conceptualize these separately as they can be independently varied—
 26 indeed, computational systems with the same affect-derivation model often have quite different
 27 intensity equations (Gratch *et al.* 2009). Intensity models usually utilize a subset of appraisal
 28 variables (e.g. most intensity equations involve some notion of desirability and likelihood); how-
 29 ever, they may involve several variables unrelated to appraisal (e.g. Elliott and Siegle 1993).
 30 Although less studied than appraisal-derivation models, some research has investigated which
 31 intensity model best conforms to human behaviour (Mellers *et al.* 1997; Gratch *et al.* 2009;
 32 Reisenzein 2009).

33 Emotion/affect

34 Affect is a representation of the agent’s current emotional state. This could be a discrete emotion
 35 label, a set of discrete emotions, core affect (i.e. a continuous dimensional space), or even some
 36 combination of these factors. An important consideration in representing affect, particularly for
 37 systems that model the consequences of emotions, is whether this data structure preserves the
 38 link between appraisal factors and emotional state. As noted above in the discussion of appraisal
 39 and dimensional theories, emotions are often viewed as being about something (e.g. I am angry

.....
 2 Smith and Kirby (2009) use the term *structural model* to refer to this mapping, drawing analogy to
 structural equation modelling, the statistical technique for estimating the causal relationships between
 variables that appraisal theorists often use to derive these mappings (see Kline 2005). As the term ‘struc-
 tural model’ is often used to contrast with ‘process models’ (a distinction we ourselves use later), we prefer
 the different term.

1 at *Valerie*). Agents that model affect as some aggregate dimensional space must either preserve the
 2 connection between affect and domain objects that initiated changes to the dimensional space, or
 3 they must provide some attribution process that *post hoc* recovers a (possibly incorrect) domain
 4 object to apply the emotional response to. For example, EM (Neal Reilly 1996) has a dimensional
 5 representation of core affect (valence and arousal) but also maintains a hierarchical data structure
 6 that preserves the linkages through each step of the appraisal process to the multiple instances of
 7 discrete emotion that underlie its dimensional calculus. In contrast, WASABI (Becker-Asano and
 8 Wachsmuth 2008) breaks this link. Some models propose some hybrid. For example, EMA main-
 9 tains discrete appraisal frames that represent specific discrete emotion instances but then allow a
 10 general dimensional ‘mood’ to moderate which discrete emotion raises to the level of awareness.

11 Affect-consequent model

12 An affect-consequent model maps affect (or its antecedents) into some behavioural or cognitive
 13 change. Consequent models can be usefully described in terms of two dimensions, one distin-
 14 guishing if the consequence is inner or outer directed (cognitive versus behavioural) and the
 15 other describing whether or not the consequence feeds into a cycle (i.e. is closed-loop).

16 Emotion can be directed outward into the environment or inward, shaping a person’s thoughts.
 17 Reflecting this, behaviour consequent models summarize how affect alters an agent’s observable
 18 physical behaviour and cognitive consequent models determine how affect alters the nature or
 19 content of cognitive processes. Most embodied computational systems model the mapping
 20 between affect and physical display, such as facial expressions. For example, WASABI maps
 21 regions of core affect into one of seven possible facial expressions (Becker-Asano 2008, p. 85) and
 22 ParleE (Bui 2004) maps from an emotion state vector (the intensity of six discrete emotion labels)
 23 to a facial muscle contraction vector (specifying the motion of 19 facial action units). Emotions
 24 can also trigger physical actions. For example, in a process called problem-focused coping, EMA
 25 (Gratch and Marsella 2004a; Marsella and Gratch 2009) attempts to mitigate negative emotions
 26 by changing features in the environment that led to the initial undesirable appraisal. In contrast,
 27 cognitive consequent models change some aspect of cognition as a result of affective state. This
 28 can involve changes in how cognition processes information; for example, Gmytrasiewicz and
 29 Lisetti (2000) propose a model that changes the depth of forward projection as a function of emo-
 30 tional state in order to model some of the claimed effects of emotion on human decision-making.
 31 Cognitive consequent models can also change the content of cognitive processes; for example,
 32 EMA implements a set of emotion-focused coping strategies, such as wishful thinking, distancing,
 33 and resignation, that alter an agent’s beliefs, desires, and intentions, respectively.

34 We can further distinguish consequences by whether or not they form a cycle by altering the
 35 circumstances that triggered the original affective response. For example, a robot that merely
 36 expresses fear when its battery is expiring does not address the underlying causes of the fear,
 37 whereas one that translates this fear into an action tendency (e.g. seek power) is attempting to
 38 address the underlying cause. In this sense, both behavioural and cognitive consequences can be
 39 classified as *closed-loop* if they directly act on the emotion eliciting circumstances or *open-loop* if
 40 they do not.

41 ♦ Open-loop models are best seen as making an indirect or mediated impact on the agent’s
 42 emotional state. For example, open-looped behavioural consequences such as emotional dis-
 43 plays make sense in multi-agent setting where the display is presumed to recruit resources
 44 from other agents. For example, building a robot that expresses fear makes sense if there is a
 45 human around that can recognize this display and plug it in. Gmytrasiewicz and Lisetti’s
 46 (2000) work on changing the depth of decision-making can similarly be seen as having
 47 an indirect on emotion: by altering the nature of processing to one best suited to a certain

emotional state, the cognitive architecture is presumably in a better position to solve problems that tend to arise when in that state.

◆ Closed-loop models attempt to realize a direct impact to regulate emotion and suggest ways to enhance the autonomy of intelligent agents. Closed-loop models require reasoning about the cognitive and environmental antecedents of an emotion so that these factors can ultimately be altered. For example, EMA implements problem-focused coping as a closed-loop behavioural strategy that selects actions that address threats to goal achievement, and implements emotion-focused coping as a closed-loop cognitive strategy that alters mental state (e.g. abandons a goal) in response to similar threats. Closed-loop models naturally implement a view of emotion as a continuous cycle of appraisal, response, and re-appraisal. In EMA, an agent might perceive another agent's actions to be a threat to its goals, resulting in anger, which triggers a coping strategy that results in the goal being abandoned, which in turn lowers the agent's appraised sense of control, resulting in sadness (see Marsella and Gratch 2009).

The component model in Figure 1.2.2 is, of course, only one of many possible ways to dissect and link emotion components but we have found it pragmatically useful in our own understanding of different computation approaches, as we illustrate below. Additionally, many of the components we identify have previously been highlighted as important distinctions with the literature on emotion research. For example, Smith and Kirby (2009) highlight appraisal-derivation as an important but understudied aspect of appraisal theory. In their work they propose the term *relational model* to refer to this component, building on Lazarus's idea that appraisal is a relational construct relating the person and the environment. As 'relation' and 'relational' often has a very different meaning within computer science, we prefer a different term, 'appraisal-derivation model'. Appraisal-derivation models are frequently identified within the appraisal theory under the term *structural model*, drawing analogy to structural equation modelling (Kline 2005), the statistical technique for estimating the causal relationships between variables that appraisal theorists often use to derive these mappings. As the term 'structural model' is not emotion-specific, and is often used to contrast with 'process models' (a distinction we ourselves use later), we prefer the term appraisal-derivation model. The idea of closed-loop models has been proposed by a variety of appraisal theorists. Most recently, Brian Parkinson (2009) has used the term *transactional model* to highlight the incremental unfolding nature of emotional reactions, although we prefer the term closed-loop, again drawing on computer metaphors.

Processing assumptions

Computational appraisal models can vary, not only by which subcomponents they choose to implement, but how these individual components are realized and interact. Some computational systems make strong commitments to *how* information is processed (e.g. in parallel or sequential?). Others make strong commitments concerning *what* information is processed (e.g. states, goals, plans). Here we introduce terminology that characterizes these different processing choices.

Process specificity

Computational models vary considerably in term of the claims they make about how information is processed. At the most abstract level, a *structural model* specifies a mapping between inputs and outputs but makes no commitment to how this mapping is realized—the term structural comes from structural equation modelling (Kline 2005), a statistical method whereby the relationship between input and output can be inferred. In contrast, a process model posits specific constraints on how inputs are transformed into outputs. For example, Ortony, Clore, and Collins present a structural affect-derivation model that maps from appraisal values to an emotion label, whereas

1 Scherer's sequential checking theory is a process appraisal-derivation model that not only
 2 specifies the structure of appraisal but proposes a set of temporal processing constraints on how
 3 appraisal variables should be derived (e.g. goal relevance should be derived before normative
 4 significance). The distinction between structural and process is not clear-cut and is best seen as a
 5 continuum. Psychological process theories only specify processes to some level of detail and dif-
 6 ferent theories vary considerably in terms of their specificity. In contrast, a computational model
 7 must be specified in sufficient detail for it to be realized as working software; however many of
 8 these process details are pragmatic and do not correspond to strong theoretical claims about how
 9 such processes should be realized. For example, Elliott's affective reasoner implements affect-
 10 derivation via a set of *ad hoc* rules, but this should not be seen as a claim that appraisal should be
 11 implemented in this manner, but rather as a short cut necessary to create a working system.

12 Processing constraints can be embedded within an individual appraisal component or can
 13 emerge from the interactions of individual components. For example, Scherer's sequential check-
 14 ing theory posits temporal ordering constraints with its model of appraisal derivation. In con-
 15 trast, Gratch and Marsella's EMA model posits that emotion arises from a continuous cycle of
 16 appraisal, coping, and reappraisal and that such temporal properties arise from the incremental
 17 evolution of the person-environment relationship (for an in-depth discussion of this point see
 18 Marsella and Gratch 2009).

19 Processing constraints can be asserted for a variety of reasons. In psychology, process models
 20 are typically used to assert theoretical claims about the nature of human mental processes, such
 21 as whether appraisal is a sequential or parallel process. Within computational systems, the story
 22 is more complex. For computational systems that model human psychological processes, the aim
 23 is the same: faithfully reflect these theoretical claims into computational algorithms. For example,
 24 Marinier (2008) translates Scherer's processing assumptions into architectural constraints on
 25 how information is processed in his PEACTIDM model. However, processing constraints can be
 26 introduced for a variety of other reasons having nothing to do with fidelity to human psycho-
 27 logical findings. These include, for example, formalizing abstract mappings into precise language
 28 (Meyer 2006; Lorini and Castelfranchi 2007), proving that a mapping is computable, illustrating
 29 efficient or robust algorithms to achieve a mapping, etc.

30 Representational specificity

31 Regardless of how component models process information, computational systems vary consid-
 32 erably in the level-of-detail of the information they process. Some model emotional processes as
 33 abstract black boxes (exploring, for example, the implications of different relationships between
 34 components), whereas others get down the nitty-gritty of realizing these processes in the context
 35 of concrete application domains. This variance is perhaps easiest to see when it comes to appraisal
 36 derivation. For example, all appraisal models decompose the appraisal process into a set of indi-
 37 vidual appraisal checks. However, some models stop at this level, treating each check as a repre-
 38 sentational primitive (e.g. Sander *et al.* 2005), whereas others further decompose appraisal checks
 39 into the representational details (e.g. domain propositions, actions, and the causal relationships
 40 between them) that are necessary for an agent to appraise its relationship to the environment (e.g.
 41 Neal Reilly 1996; El Nasr *et al.* 2000; Gratch and Marsella 2004a; Dias and Paiva 2005; Mao and
 42 Gratch 2006; Becker-Asano 2008; Si *et al.* 2008).

43 Process specificity can vary independently from representational specificity. For example,
 44 Sander and colleagues (2005) provide a detailed neural network model of how appraisals
 45 are derived from the person-environment relationship, but the person-environment relationship
 46 itself is only abstractly represented. Process and representational specificity also vary across com-
 47 ponent models within the same system. For example, WASABI (Becker-Asano 2008)

1 incorporates detailed representational and process commitments for its model of affect-
 2 derivation but uses less detail for its model of appraisal derivation. Such differences often result
 3 from the fact that, while specific systems are directed at addressing a subset of the components
 4 involved in emotion processes, the authors often require a complete working system to assess
 5 the impact of their proposed contribution and these other components may be rudimentary or
 6 *ad hoc*.

7 Domain specific versus domain independent

8 In addition to their processing constraints and representational specificity, algorithms can be
 9 characterized by the generality of their implementation. A domain-independent algorithm
 10 enforces a strict separation between details of a specific domain, typically encoded as a *domain*
 11 *theory*, and the remaining code, which is written in such a way that it can be used without modi-
 12 fication. For example, planning algorithms operate on a domain theory consisting of a set of
 13 states and actions that describe a domain and provide general algorithms that operate syntacti-
 14 cally on those representations to generate plans. Computational appraisal models differ in terms
 15 of how domain-specific knowledge is encoded and which components require domain-specific
 16 input. Most systems incorporate domain-independent affect-derivation models (Neal Reilly
 17 1996; Bui 2004; Gratch and Marsella 2004; Gebhard 2005; Becker-Asano 2008; Marinier 2008).
 18 Fewer systems provide domain-independent algorithms for appraisal-derivation (e.g. Neal Reilly
 19 1996; El Nasr *et al.* 2000; Gratch and Marsella 2004; Si *et al.* 2008).

20 Example of applications of this framework

21 Viewing a computational model of emotion as a model of models allows more meaningful
 22 comparisons between systems. Systems that appear quite different on the surface can be seen
 23 as adopting similar choices along some dimensions and differing in others. Adopting a compo-
 24 nent model framework can help highlight these similarities and differences, and facilitate empiri-
 25 cal comparisons that assess the capabilities or validity of alternative algorithms for realizing
 26 component models.

27 Table 1.2.1 illustrates how the component model framework can highlight conceptual
 28 similarities and differences between emotion models. This table characterizes three quite different
 29 systems. EMA is the authors' own work on developing a general computational model of appraisal
 30 and coping motivated by the appraisal theory of Richard Lazarus (Lazarus 1991) and has been
 31 applied to driving the behaviour of embodied conversational agents (Swartout *et al.* 2006);
 32 FLAME is an OCC-inspired appraisal model that drives the behaviour of characters in interactive
 33 narrative environments (El Nasr *et al.* 2000); and ALMA is intended as a general programming
 34 tool to allow application developers to more easily construct computational models of emotion
 35 for a variety of applications (Gebhard 2005). Some observations that can be made from this table
 36 include the following.

- 37 ♦ EMA and FLAME both focus on appraisal derivation. They provide domain-independent
 38 techniques for representing the person–environment relationship and derive appraisal varia-
 39 bles via domain-independent inference rules, although the approaches adopt somewhat dif-
 40 ferent representational and inferential choices. In contrast, ALMA does not address appraisal
 41 derivation.
- 42 ♦ All systems in Table 1.2.1 use rules to derive affect from a set of appraisal variables. ALMA and
 43 FLAME both adopt OCC-style appraisal variables and affect-derivation models, whereas EMA
 44 uses a model inspired by Lazarus.
- 45 ♦ Each system adopts a different choice for how the intensity of an emotion is calculated.

Table 1.2.1 A component model view of three different systems

	EMA	ALMA	FLAME
Person–environment relationship	Domain-independent decision-theoretic plans + BDI	Outside the scope of model	Domain-independent Markov-decision process
Appraisal-derivation	Inference over decision-theoretic plans	User-defined	Fuzzy rules over Markov-decision graph
Appraisal-variables	Lazarus-Inspired, desirability, likelihood, expectedness, causal attribution, controllability, changeability	OCC-inspired Good/bad, likely/unlikely event Good/bad act of self/other Nice vs. nasty thing	OCC-inspired Desirability Expectation (dis)approval
Affect-derivation	Lazarus-based structural model that generates discrete emotion and mood state	OCC-based structural model that gives 'Impulsed' into core affect	OCC-based structural model producing discrete emotion labels
Affect-intensity	Expected utility model, Threshold model, Additive mood derivation	User defined	Additive model
Affect	Set of appraisal frames, mood (discrete-emotion vector) with decay	PAD space representing both current mood and emotion	Emotion and positive vs. negative mood state
Behavioural-consequences	Most-intense emotion alters behavior display and action selection. Actions are close-loop via domain-independent rules	Open looped. Mood and emotion alter behaviour display and action selection	Domain-specific fuzzy expression and action rules
Cognitive-consequences	Closed-loop via domain-independent emotion-focused coping that changes BDI	Open-looped. Emotion amplifies/dampens intensity of elicited emotions	Closed-loop changes to domain model via reinforcement learning

- 1 ♦ All systems incorporate some notion of core affect, though they adopt different representations. EMA has a mood state that summarizes the intensity of all active emotional appraisals
- 2 and this mood biases the selection of a single emotional appraisal that can impact behaviour.
- 3 ALMA represents both a current emotion and a more general mood in a three-dimensional
- 4 (PAD) space (either of which can impact behaviour). FLAME has a one-dimensional (positive
- 5 vs. negative) mood state that can influence behaviour.
- 6
- 7 ♦ EMA and FLAME propose closed-loop consequence models that allow emotion to feed back
- 8 into changes in the mental representation of a situation, although they adopt quite different
- 9 algorithmic choices for how to realize this function.
- 10 Besides allowing such conceptual comparisons, the key benefit of decomposing systems into
- 11 component models is that it allows individual design decisions to be empirically assessed inde-
- 12 pendent of other aspects of the system. For example, in Table 1.2.1, FLAME and EMA adopt dif-
- 13 ferent models for deriving the intensity of an emotional response: both systems calculate intensity
- 14 as a function of the utility and probability of goal attainment but FLAME adds these variables

Table 1.2.2 Comparing the fit of different affect intensity models

	Hope	Joy	Fear	Sadness
Expectation-Change model	PEACTIDM	ParleE, EM PEACTIDM	PEACTIDM	ParleE, EM PEACTIDM
Expected utility	EMA, ParleE, FearNot! EM BTDE		EMA, ParleE, FearNot! EM BTDE	
Threshold model		EMA, FearNot! BTDE		EMA, FearNot!
Additive model	Cathexis, FLAME	Cathexis, FLAME	Cathexis, FLAME	Cathexis, FLAME
Hybrid model	Price <i>et al.</i> , 1985	Price <i>et al.</i> , 1985	Price <i>et al.</i> , 1985	Price <i>et al.</i> , 1985

1 whereas EMA multiples them for prospective emotions (e.g. hope and fear) and uses a threshold
 2 model for retrospective emotions (e.g. joy and sadness). An advantage of the component model
 3 view is these alternative choices can be directly compared and evaluated, independently of the
 4 other choices adopted in the systems from which they stem.

5 Gratch and Marsella recently applied this component-model perspective to an empirical
 6 comparison of different affect-derivation models (Gratch *et al.* 2009). Besides the two approaches
 7 proposed by EMA and FLAME, researchers have proposed a wide range of techniques to calculate
 8 the intensity of an affective response. In their study, Gratch and Marsella analysed several com-
 9 peting approaches for calculating the intensity of a specific emotional response to a situation and
 10 classified these approaches into a small number of general approaches (this includes approaches
 11 used in a variety of systems including: Price *et al.* 1985; Neal Reilly 1996; Velásquez 1998; El Nasr
 12 *et al.* 2000; Bui 2004; Dias and Paiva 2005; Marinier *et al.* 2009; Reisenzein 2009). They then
 13 devised a study to empirically test these competing appraisal intensity models, assessing their
 14 consistency with the behaviour of a large number of human subjects in naturalistic emotion-
 15 eliciting situations. In the study they had subjects play a board game (Battleship™ by Milton
 16 Bradley™) and assessed subjects' self-reported emotional reactions as the game unfolded and as a
 17 consequence of if they were winning or losing (which was manipulated experimentally).

18 Table 1.2.2 summarizes the results of this study that compared the behaviour of EMA to several
 19 other systems proposed in the literature. These include ParleE (Bui 2004), a system that uses
 20 appraisal models to drive facial animation; BTDE (Reisenzein 2009), an appraisal theory that
 21 attempts to reduce appraisal-derivation, affect-derivation, and affect-intensity to operations over
 22 beliefs and desires; FLAME, described above; Cathexis, an anatomical approach that views emo-
 23 tions as arising from drives; FearNot! (Dias and Paiva 2005), an appraisal model based on EMA;
 24 EM (Neal Reilly 1996) an OCC-inspired model that drives the behaviour of interactive game
 25 characters; and a model proposed by Price and colleagues (Price, Barrell, *et al.* 1985) that inspired
 26 the design of FLAME. Although these models vary in many ways, when it comes to affect-intensity,
 27 they can be described in terms of four basic methods have been proposed in the literature for
 28 deriving the intensity of an emotional response.

29 As noted in the table, different systems used different intensity models depending on the
 30 emotion type. The intensity models, listed in the first column, include expected utility (i.e. the
 31 intensity of emotional response is proportional to the utility of a goal times its probability of
 32 attainment), expectation-change (i.e. the intensity is proportional to the change in probability
 33 caused by some event), and additive (i.e. the intensity is proportional to the sum of probability
 34 and utility). The cells in the table indicate the intensity model that a particular system applies to

1 calculate the intensity of a given emotion. The table also summarizes the results of how well these
 2 different models explain the data elicited from the study. A slash through the box indicates the
 3 model cannot explain the results of the experiment. This analysis lends support for the expected
 4 utility model for all emotions, with a particularly strong fit for the prospective emotions (i.e. hope
 5 and fear), though it allows that some modified version of a threshold model might explain the
 6 results of retrospective motions like joy and sadness. If the goal of an emotion model is to realisti-
 7 cally model human emotional responses, expected utility is probably a good choice for that
 8 appraisal intensity component model.

9 Of course, the behaviour of a specific component is not necessarily independent of other design
 10 choices so such a strong independence assumption should be treated as a first approximation for
 11 assessing how alternative design choices will function in a specific system. However, unless there
 12 is a compelling reason to believe choices are correlated, such an analysis should be encouraged.
 13 Indeed, a key advantage of the compositional approach is that it forces researchers to explicitly
 14 articulate what these dependencies might be, should they wish to argue for a component that is
 15 repudiated by an empirical test that adopts a strong assumption of independence.

16 Dividing computational emotion models into components enables a range of such empirical
 17 studies that can assess the impact of these design choices on the possible uses of emotion models
 18 that were outlined at the start of this chapter—i.e. their impact on psychological theories of emotion;
 19 their impact on artificial intelligence and robotics; and their impact on human–computer inter-
 20 action. Here we touched on some studies that more naturally apply to the first goal and several
 21 examples of this now exist including evaluations of the psychological validity of cognitive conse-
 22 quent models (Marsella *et al.* 2009) and appraisal-derivation models (e.g. Mao and Gratch 2006;
 23 Tomai and Forbus 2007). However, the same approach can be equally applied to these other
 24 overall objectives. For example, de Melo and colleagues present evidence that the appraised
 25 expression of emotion can influence human–computer interaction in the context of social games
 26 such as iterated prisoner’s dilemma (de Melo *et al.* 2009) and it would be interesting to consider
 27 how different appraisal-derivation and intensity models might impact the power of this effect.
 28 Other researchers have explored how emotions might improve the decision-making capabilities
 29 of general models of intelligence (Scheutz and Sloman 2001; Ito *et al.* 2008) and a component
 30 model analysis can give greater insight into which aspects of these models contribute to enhanced
 31 performance.

32 Conclusion

33 In this chapter, we provided an overview of research into computational models of emotion that
 34 details the common uses of the models and the underlying techniques and assumptions from
 35 which the models are built. Our goals were twofold. For researchers outside the field of computa-
 36 tional models on emotion, we want to facilitate an understanding of the field. For research in the
 37 field, our goal is to provide a framework that can help foster incremental research, with researchers
 38 relying on careful comparisons, evaluations, and leveraging to build on prior work, as a key to
 39 forward progress.

40 To achieve those goals, we presented several conceptual distinctions that can aid in evaluation
 41 of competing models. We identified several roles for models, in psychological research, in human–
 42 computer interaction, and in AI. Evaluation, of course, must be sensitive to these roles. If, for
 43 example, the model is being used as a methodological tool for research in human emotions or in
 44 human–computer interaction research as a means to infer user emotional state, then fidelity of the
 45 model with respect to human behaviour will be critical. If the model is to be used to create virtual



1 characters that facilitate engagement with, or influence of, humans then fidelity may be less
2 important, even undesirable, while effectiveness in the application becomes more important.
3 Our assumption is that, regardless of the specific details of the evaluation, research progress in
4 computational models of emotion critically hinges not only on evaluations of specific models but
5 also on the comparison across models. Due to the complexity of some of these models, and their
6 emphasis on different aspects of the overall emotion process, it may not be reasonable or desirable
7 to undertake comparison and evaluation *in toto*. Rather component-by-component analyses,
8 based on a common lexicon, will be both more revealing and often easier. Our hope is that the
9 application of the component analyses exemplified above may serve as a means to facilitate this
10 component-by-component evaluation and lead to additional work in this direction.



