# GENERALIZED NON-ABELIAN RECIPROCITY LAWS: A CONTEXT FOR WILES' PROOF

## AVNER ASH AND ROBERT GROSS

#### Abstract

In as elementary a way as possible, we place Wiles' proof of Fermat's last theorem into the context of a general description of reciprocity conjectured to obtain between algebraic varieties defined over  $\mathbf{Q}$  and Hecke eigenvectors in the homology of the spaces of lattices in  $\mathbf{R}^n$ .

We shall find the Cube of the Rainbow, Of that, there is no doubt. But the Arc of a Lover's conjecture Eludes the finding out.

#### Emily Dickinson

During the last few decades, the field of number theory has been increasingly permeated by the theory of automorphic forms and automorphic representations. This phenomenon often goes under the rubric of the 'Langlands program', although it involves the work of many mathematicians, including R. Langlands, J.-P. Serre and G. Shimura, to name only three. Recently, this program became especially prominent because it forms a background to Wiles' proof of Fermat's last theorem.

As explained by Langlands in [14], parts of this program can be viewed as a vast generalization of reciprocity laws familiar in number theory, such as quadratic and Artin reciprocity. Some of the ultimate conjectures along these lines are spelled out in Clozel's article [4], linking *motives* and *automorphic representations*; see also Gelbart [11].

Both of the italicized terms in the previous sentence are rather technical objects. Our goal in this article is to describe a version of these generalized non-abelian reciprocity conjectures: a version accessible to a reader who knows nothing beyond basic algebra and the definition of homology groups. Because of these self-imposed limitations, we shall be unable to state the conjectures in their full strength or with total precision. To compensate for these necessary shortcomings, we have included three basic examples that should give the flavour of the conjectures. The first involves quadratic reciprocity, the second a modular elliptic curve, and the third a (probably) automorphic algebraic surface. Moreover, the conjecture that we do state is strong enough to allow its application to Fermat's last theorem. At the appropriate place, we shall discuss the term 'reciprocity' and what may be thought of as being reciprocated.

By reducing the prerequisites, we hope to increase the number of readers who can obtain some glimpses of this beautiful landscape of generalized reciprocity. We also believe that the formulation of the conjectures in terms of spaces of lattices

Bull. London Math. Soc. 32 (2000) 385-397

Received 30 July 1997; revised 29 November 1999.

<sup>2000</sup> Mathematics Subject Classification 11-02.

in  $\mathbf{R}^n$  has a certain surprising beauty of its own. We should point out, however, that it is very unlikely that progress can be made in *proving* the conjectures on this elementary level.

We begin with a quick review of Wiles' proof of Fermat's last theorem. For a good survey of this together with references to the work of Wiles and his predecessors necessary to the proof, see [15].

Let *n* be a prime greater than 3, and suppose that  $(a, b, c) \in \mathbb{Z}^3$  is a non-trivial solution of the Fermat equation  $a^n + b^n = c^n$ . It was noticed that one could take these three integers and use them to define a particular elliptic curve *E*,

$$y^2 = x(x - a^n)(x + b^n),$$

where we think of *E* as lying in the *xy*-plane.

Elliptic curves have been the object of intense study for much of this century, though many things about them remain unknown. One particular property of elliptic curves is that they can be *modular*. We shall give below a definition of modularity from our point of view that is equivalent to the standard one. Wiles proved that this particular elliptic curve is modular; Ribet had earlier proved that E is not modular. Therefore it cannot exist, so a, b and c cannot exist either.

The work of Wiles [23], Taylor and Wiles [21], Diamond [7], and Conrad, Diamond and Taylor [5] shows that elliptic curves in a rather large class are modular. Just recently, Christophe Breuil, Brian Conrad, Fred Diamond and Richard Taylor have announced a proof of the modularity conjecture [6], which asserts that every elliptic curve defined over Q is modular.

This conjecture is part of a vast program which sets up a conjectural correspondence between two large sets. On one side of the correspondence, roughly speaking, is the collection of systems of simultaneous polynomial equations with rational coefficients. A member of this collection is called a 'variety defined over  $\mathbf{Q}$ ', and is relatively simple to comprehend. (More precisely, this side of the correspondence contains motives, which we shall discuss below.) For example, elliptic curves defined over  $\mathbf{Q}$  are varieties, because they can be described in general by an equation

$$y^{2} + a_{1}xy + a_{3}y = x^{3} + a_{2}x^{2} + a_{4}x + a_{6}$$

where the  $a_i$  are integers, and the discriminant  $\Delta$  of the elliptic curve, which is a polynomial function of the numbers  $a_i$ , is non-zero. Points on the curve over the ring R are simply solutions  $(x, y) \in R^2$  to the given equation.

The formula for the discriminant is quite complex, and may be found in [19, p. 46]; if we restrict to the simpler family of curves given by the equation

$$y^2 = x^3 + Ax + B,$$

then  $\Delta = -16(4A^3 + 27B^2)$ . Except for a constant factor, this is the same as the discriminant of the cubic polynomial  $x^3 + Ax + B$ . For future reference, we mention a more subtle number N called the *conductor*, which consists of a product of the primes dividing  $\Delta$  raised to certain powers.

The other side of the correspondence is harder to grasp; it consists of certain automorphic representations of reductive groups defined over  $\mathbf{Q}$ . (Linking these two collections is something even harder to get an elementary handle on: representations of Gal( $\overline{\mathbf{Q}}/\mathbf{Q}$ ), which we shall discuss further in the Appendix.) To simplify our exposition, in this article we shall consider only a certain subset of automorphic representations: those which are 'geometric' for GL(*n*) with constant coefficients. This

particular subset is still rich enough to allow us to define the concept of modularity for elliptic curves. This narrowing of scope allows us to replace automorphic representations by simpler equivalent objects, called **H**-eigenelements, to be explained below. We denote by  $\alpha$  an element of this collection.

The first side of the correspondence will now consist of a more narrowly defined collection of varieties that will still include elliptic curves. We shall denote such a variety by V. Our subprogram of the general reciprocity program is the following correspondence:

$$\begin{pmatrix} \text{certain varieties (including} \\ \text{elliptic curves) defined over } \mathbf{Q} \end{pmatrix} \iff (\mathbf{H}\text{-eigenelements}) \\ V \iff \alpha.$$

We need to explain more precisely the terms and nature of this correspondence. First, the left-hand side: by multiplying through by the common denominator of all fractions used to define the variety V, we may assume that in fact the equations used to define V have coefficients in  $\mathbb{Z}$  rather than in  $\mathbb{Q}$ . If R is any ring, we can then use V(R) to mean the solutions to the equations which are in R. In other words, if V is defined by polynomials  $\{g_i(x_1, \ldots, x_m) = 0, i = 1, \ldots, n\}$ , then

$$V(R) = \{(a_1, \dots, a_m) \in R^m \mid g_i(a_1, \dots, a_m) = 0 \text{ for } i = 1, \dots, n\}.$$

In particular, we can take R to be  $\mathbf{F}_{p^m}$ , the field of  $p^m$  elements, where p is any prime and m is any positive integer. Because the field is finite, we obtain a set of integers  $\#V(\mathbf{F}_{p^m})$  for primes p and positive integers m. These integers are encoded in the Hasse–Weil zeta function of V as described below.

Unfortunately, making precise which equivalence classes of varieties are considered on the left-hand side of this correspondence is quite complex. In fact, we actually need 'pieces' of varieties, or, more precisely, 'pieces' of their cohomology, called *motives*. When the cohomology of a variety V breaks up into a number of motives, all of which except one are easily understood, then V as a whole can stand in for its non-trivial motivic piece in the formulas that give a description of the correspondence  $V \leftrightarrow \alpha$ . This is what happens in the three examples discussed in this paper. Therefore, the first-time reader might wish to skim this *ad hoc* introduction to motives. Advanced readers may consult [12] for more information about motives.

A motive is a piece of the cohomology of a variety defined over  $\mathbf{Q}$ . Just like the complex cohomology of a complex variety, motives have Hodge types, to which we shall need to refer for the sake of accurately stating the conjectures below. One projective variety V can 'contain' many motives, and one motive can appear in various different varieties. A motive M gives rise to an L-function L(M,t) with an Euler product

$$L(M,t) = L_{\infty}(M,t) \prod_{p} L_{p}(M,t),$$

where p runs over all prime numbers. If the variety V yields the set of motives  $\{M_i\}$ , then

$$\prod_{i} L(M_i, t)^{\varepsilon_i} = Z(V, t) = Z_{\infty}(V, t) \prod_{p} Z_p(V, t),$$

where Z(V, t) is the Hasse-Weil zeta function of V, and  $\varepsilon_i$  is -1 or 1 depending on whether  $M_i$  appears in an odd- or even-degree cohomology group of V. The local

factor  $Z_p(V, t)$  is defined by

$$Z_p(V,t) = \exp\left(\sum_{m=1}^{\infty} \frac{\#V(\mathbf{F}_{p^m})}{m} t^m\right).$$

In fact, we even have the local factorization

$$Z_p(V,t) = \prod_i L_p(M_i,t)^{\varepsilon_i}$$

for each prime p.

An example could not but help to clarify these concepts. Consider the variety  $V = \mathbf{P}^2$ , defined over  $\mathbf{Q}$ . We can think of  $\mathbf{P}^2$  as the disjoint union  $\mathbf{A}^0 \amalg \mathbf{A}^1 \amalg \mathbf{A}^2$ . Therefore  $\#\mathbf{P}^2(\mathbf{F}_{p^m}) = 1 + p^m + p^{2m}$ . We have

$$\log Z_p(V,t) = \sum_{m=1}^{\infty} \frac{1+p^m+p^{2m}}{m} t^m$$
$$= \sum_{m=1}^{\infty} \left[ \frac{t^m}{m} + \frac{(pt)^m}{m} + \frac{(p^2t)^m}{m} \right]$$
$$= -\left( \log(1-t) + \log(1-pt) + \log(1-p^2t) \right),$$

so

$$Z_p(V,t) = \frac{1}{(1-t)(1-pt)(1-p^2t)}.$$

It is traditional to substitute  $t = p^{-s}$ , and view s as a complex variable. We then have

$$Z_p(V, p^{-s}) = \left( (1 - p^{-s})(1 - p^{1-s})(1 - p^{2-s}) \right)^{-1},$$

and therefore

$$\prod_{p} Z_p(V, p^{-s}) = \zeta(s)\,\zeta(s-1)\,\zeta(s-2).$$

In this case, the three motives corresponding to V are exactly  $H^i(V)$  for i = 0, 2, 4and  $L(M_i, p^{-s}) = \zeta(s - i/2)$ .

In general, for any smooth projective variety V defined by equations with integral coefficients, and for almost all primes p,  $Z_p(V, t)$  is a rational function with

$$Z_p(V,t) = \frac{P_{1,p}(t) P_{3,p}(t) \cdots}{P_{0,p}(t) P_{2,p}(t) \cdots}$$

with the factorization

$$P_{i,p}(t) = \prod_{r} (1 - \beta_{r,i}^{(p)} t).$$

Let  $\overline{V} = V(\overline{\mathbf{F}_p})$ . Then the Frobenius element  $\phi_p$  acts on the étale cohomology of  $\overline{V}$ , with  $\beta_{r,i}^{(p)}$  the reciprocal eigenvalues of  $\phi_p$ . We have  $|\beta_{r,i}^{(p)}| = p^{i/2}$ . Finally,

$$\#V(\mathbf{F}_{p^m}) = \sum_{r,i} (-1)^i (\beta_{r,i}^{(p)})^m$$

Now, a motive *M* corresponds first to a choice of a single index *i* and a positive integer *n*, and then for each prime *p* to a choice of a subset  $S_p$  of  $\{\beta_{r,i}^{(p)}\}$  of cardinality

*n* such that  $S_p$  is the set of reciprocal eigenvalues of  $\phi_p$  acting on a 'geometrically defined piece' of the étale cohomology of *V*. For example, the 'piece' might be all of  $H^i(V)$ . See the discussion of [10] later in this paper for an example of a motive that is not all of  $H^i(V)$ .

We can now define the L-function of a motive M by setting

$$L_p(M,s) = \prod_{\beta \in S_p} (1 - \beta p^{-s})^{-1}$$
, for almost all  $p$ .

Note the customary change of variable from t to s. We omit the definitions for 'bad' primes and for  $L_{\infty}$ , but we mention that  $L_{\infty}(M, s)$  depends on the Hodge type of M. Then

$$L(M,s) = L_{\infty}(M,s) \prod_{p} L_{p}(M,s).$$

Now we can say that the left-hand side of the correspondence consists of motives of dimension *n* with Hodge type  $(n-1,0) \oplus (n-2,1) \oplus \cdots \oplus (0,n-1)$ .

On the right-hand side of the correspondence, **H** will denote a ring of operators acting on certain homology groups, and this action has eigenelements. The eigenvalues of the operators in **H** acting on  $\alpha$  will also be a set of numbers, and the reciprocity conjecture specifies a relationship between these eigenvalues and the numbers  $\#V(\mathbf{F}_{p^m})$ . These relationships are given precisely by an equality of *L*-functions; but first we have to define the  $\alpha$  and their *L*-functions.

What are these mysterious H-eigenelements? We begin with a review of a common definition.

DEFINITION. A *lattice*  $\Lambda \subset \mathbf{R}^n$  is a free abelian group generated by a basis of  $\mathbf{R}^n$ .

Our next definition is less common, but also simple to comprehend.

DEFINITION. A level N structure on an n-dimensional lattice  $\Lambda$  is a group isomorphism  $\phi : \Lambda/N\Lambda \to (\mathbb{Z}/N\mathbb{Z})^n$ .

We could consider the collection of all lattices in  $\mathbb{R}^n$  with level N structure, but that turns out to be too large a class. Instead, we consider two such lattices to be equivalent if one can be obtained from the other by proper Euclidean motion (that is, an orthogonal transformation of determinant 1), and positive homothety (change of scale), where we always view  $\mathbb{R}^n$  as endowed with its usual Euclidean structure.

DEFINITION.  $L_n(N) = \{(\text{lattice } \Lambda \subset \mathbf{R}^n, \text{ level } N \text{ structure } \phi)\} / \langle \text{proper Euclidean motion, positive homotheties} \rangle.$ 

Before we consider a few examples, we put a topology on our set  $L_n(N)$ . Given  $(\Lambda, \phi) \in L_n(N)$ , fix a basis for  $\Lambda$ , and then consider nearby points to be those lattices obtained by small perturbations of the basis in  $\mathbb{R}^n$ , keeping the level N structure  $\phi$  the same. These spaces of lattices come into the game because they are locally symmetric spaces on which automorphic forms naturally live.

The simplest example is  $L_1(1)$ . Because our lattices are equivalent up to homothety, we can take a basis of our lattice  $\Lambda$  to be 1. A level 1 structure would be an isomorphism from  $\Lambda/\Lambda$  to  $\mathbb{Z}/\mathbb{Z}$ , and since each group contains only the identity element, there can be only one such isomorphism. Therefore  $L_1(1)$  consists of a single point.

A more enlightening example is given by  $L_1(N)$ . Again, there is only one lattice  $\Lambda$  up to homothety, so we can take  $\Lambda = \mathbb{Z}$ . Our level N structure is an isomorphism  $\phi : \mathbb{Z}/N\mathbb{Z} \to \mathbb{Z}/N\mathbb{Z}$ . Such a map is defined by  $\phi(1)$ , and since  $\phi$  must be invertible,  $\phi(1)$  must be invertible, and hence it is in  $(\mathbb{Z}/N\mathbb{Z})^*$ . Therefore  $L_1(N) \cong (\mathbb{Z}/N\mathbb{Z})^*$ .

In order to see how  $L_n(N)$  can contain more interesting geometric information, we consider one more example:  $L_2(1)$ . We start with a two-dimensional lattice  $\Lambda$ . We take a vector of minimal length in  $\Lambda$ , and by means of rotations and stretching, we can take that vector to be (1,0). We are free to take any linearly independent vector in  $\Lambda$  as the second basis vector. Take any such vector (x, y), where we may suppose that  $x \neq 0$  and y > 0. By adding multiples of (1,0) to this vector, we may suppose that  $-\frac{1}{2} \leq x \leq \frac{1}{2}$ . Since this vector must have length at least as large as (1,0), we also know that  $x^2 + y^2 \ge 1$ .

Furthermore, there are two identifications that we can make. Clearly,  $(-\frac{1}{2}, y)$  is equivalent to  $(\frac{1}{2}, y)$ . Less obviously, there is another identification: if the second chosen basis vector (x, y) has length 1, then we can rotate this vector to (1, 0), and the vector formerly at (1, 0) will rotate to (x, -y). After we multiply by -1 to obtain a positive second coordinate, we see that (x, y) is equivalent to (-x, y) when  $x^2 + y^2 = 1$ .

After making these identifications, our space of lattices is topologically equivalent to a sphere with one point removed. In fact, taking the strip

$$\{(x, y) \in \mathbf{R}^2 \mid -\frac{1}{2} \leq x \leq \frac{1}{2}, \ x^2 + y^2 \ge 1, \ y > 0\},\$$

and identifying the left and right vertical edges, gives a cylinder. Then glueing (x, y) to (-x, y) for those points (x, y) with  $x^2 + y^2 = 1$  collapses the circle at the bottom of the cylinder to an arc. Topologically, we now have an open cup, which is homeomorphic to a sphere minus a point. (A more detailed discussion can be found in many places in the literature; see, for example, [16, Chapter VII.1] for a different explanation and a picture.) As before, the level 1 structure does not add any further detail, since  $\phi$  will be a map from the trivial group to the trivial group.

Although the general picture is obviously much more complex, certain features of this situation will remain true. For all  $n \ge 1$ ,  $N \ge 1$ ,  $L_n(N)$  will be a 'nice' topological space with  $\#(\mathbb{Z}/N\mathbb{Z})^*$  components, classified by det  $\phi$ . The situation for n = 2 is especially favourable:  $L_2(N)$  is the disjoint union of Riemann surfaces, and can be defined by equations with Q-coefficients. There are formulas for the genus of  $L_2(N)$ ; see [18, Chapter 1]. If  $n \ge 3$ , then  $L_n(N)$  is not an algebro-geometric space, but only a V-manifold (a manifold if  $N \ge 3$ ).

We next consider the homology of these spaces:  $H_d(L_n(N), \mathbb{C})$ , which is the group of *d*-dimensional cycles modulo *d*-dimensional boundaries. This is always a finite-dimensional vector space.

There is extra structure given by the action of the *Hecke algebra*  $\mathbf{H}$ , defined as follows. Let p be any prime not dividing N, and let k be any integer between 0 and n, inclusive. We can then define

$$T_p^{(k)} : L_n(N) \longrightarrow L_n(N)$$
  
( $\Lambda, \phi$ )  $\longmapsto \{ (\Lambda', \psi) \mid \Lambda' \subset \Lambda, \ \Lambda/\Lambda' \cong (\mathbf{Z}/p\mathbf{Z})^k, \ \psi = \phi|_{\Lambda'} \}.$ 

There are several observations to be made about these operators  $T_p^{(k)}$ .

- (1) Because  $p \nmid N$ ,  $\psi$  is again a level N structure.
- (2)  $T_p^{(k)}$  is not a function, but rather a one-to-many map. In fact, the image of  $(\Lambda, \phi)$  is finite, and the cardinality is given by the number of k-planes in  $(\mathbb{Z}/p\mathbb{Z})^n$ .
- (3) The action on the homology groups is, in fact, a function. In fact, T<sup>(k)</sup><sub>p</sub> maps a cycle to the union of all of its images, but that union is again a cycle. Therefore we do have a well-defined function T<sup>(k)</sup><sub>p</sub> : H<sub>d</sub>(L<sub>n</sub>(N), C) → H<sub>d</sub>(L<sub>n</sub>(N), C).
- (4) These functions  $T_p^{(k)}$  all commute. Therefore there are simultaneous eigenclasses  $\alpha$ . If we write  $T_p^{(k)}(\alpha) = a_p^{(k)}\alpha$ , then one can prove that these numbers  $a_p^{(k)}$  are algebraic integers. For a fixed  $\alpha$ , they in fact generate a finite-degree extension of **Q**.
- (5) If N|N' and  $\alpha \in H_d(L_n(N), \mathbb{C})$  is an **H**-eigenelement, then there always exists  $\alpha' \in H_d(L_n(N'), \mathbb{C})$ , an **H**-eigenelement with eigenvalues equal to those of  $\alpha$  for those primes p not dividing N'.
- (6) In principle, these eigenclasses  $\alpha$  and numbers  $a_p^{(k)}$  are computable. For n = 1 and any N, the computation is easy. For n = 2 and relatively small N, the computation is not impossible; there is much help coming from the theory of modular forms and elliptic curves. For n = 3 and relatively small N, a computer can provide the answers [1]. For n = 4 and very small N, there is some research in progress by the first author, Paul Gunnells and Mark McConnell.

In theory, whenever one can obtain a correspondence between a variety V (really a motive M) and such an  $\alpha$ , there is information to be gained. Such correspondences are 'generalized reciprocity laws', so called because quadratic reciprocity can be interpreted as such a correspondence, as we shall see shortly.

The term 'reciprocity' seems to go back to Legendre, as quoted on p. 328 of [22]. Originally, the term referred to reciprocity between two primes p and q: whether or not p was a square modulo q being determined according to a simple rule depending on whether or not q was a square modulo p. Eventually, this was interpreted as a property of  $\phi_p$ , the Frobenius element at p, restricted to the Galois group of  $\mathbf{Q}(\sqrt{\pm q})$  over  $\mathbf{Q}$ , and vice versa. Later, the term 'reciprocity' was extended to a variety of rules that told how  $\phi_p$  acted in various situations; see [24] for a good introduction. An early example of a non-abelian reciprocity law was given by Shimura [17]. Its modern use is explained by Langlands [14, pp. 408–409] as including assertions 'that an *L*-function defined by diophantine data, that is, by an algebraic variety over a number field, is equal to an *L*-function defined by analytic data, that is, by an automorphic form'. See also Tate's article [20] in the same volume, and the Appendix below, for some examples of what kind of concrete information a reciprocity law can provide.

We can now offer a precise conjecture. An **H**-eigenvector corresponds to a motive as follows. We can define an *L*-function corresponding to the **H**-eigenelement  $\alpha$  by defining

$$L_p(\alpha, s) = 1 - a_p^{(1)} p^{-s} + a_p^{(2)} p^{1-s} - a_p^{(3)} p^{3-s} + \dots + (-1)^n a_p^{(n)} p^{\frac{1}{2}n(n-1)-s}$$

for almost all p (that is, for those finite primes p not dividing the level N), with other factors for primes dividing N and for  $\infty$ , and then define

$$L(\alpha, s) = L_{\infty}(\alpha, s) \prod_{p} L_{p}(\alpha, s).$$

Conjecturally,  $\alpha$  corresponds to a motive M in the sense that  $L(\alpha, s) = L(M, s)$ . More precisely, we have the following.

CONJECTURE. (i) Given an absolutely irreducible motive M of dimension n with Hodge type  $(n-1,0) \oplus (n-2,1) \oplus \cdots \oplus (0,n-1)$  with conductor N, there is a cuspidal **H**-eigenclass  $\alpha \in H_*(L_n(N), \mathbb{C})$  such that  $L(\alpha, s) = L(M, s)$ .

(ii) Given a cuspidal **H**-eigenclass  $\alpha \in H_*(L_n(N), \mathbb{C})$ , there is a motive M of dimension n and conductor N and Hodge type  $(n - 1, 0) \oplus (n - 2, 1) \oplus \cdots \oplus (0, n - 1)$  such that  $L(\alpha, s) = L(M, s)$ .

We make some observations.

- (1) Part (i) is Question 4.16 in [4], and part (ii) is Conjecture 4.5 in [4]. We have restricted each of Clozel's formulations to the case of a trivial coefficient module for the homology of  $L_n(N)$ .
- (2) Other Hodge types occur if we allow the cohomology of  $L_n(N)$  with non-trivial coefficient modules.
- (3) The conductor N of a motive M is a well-defined attribute of M, independent of the Conjecture.
- (4) 'Cuspidal' is a technical term whose definition is beyond the scope of this paper. Roughly speaking, an H-eigenclass α ∈ H<sub>\*</sub>(L<sub>n</sub>(N), C) is cuspidal if it cannot be 'induced' from any space of lattices in R<sup>m</sup> for m < n. If n = 2 or 3, then the concept simplifies: α is non-cuspidal if for any compact subset K of L<sub>n</sub>(N), α is homologous to a cycle supported outside K. If α is cuspidal, then it is known that L(α, s) has an analytic continuation to the entire s-plane, and therefore the Conjecture implies that L(M, s) also has an analytic continuation.
- (5) The equality of *L*-functions is equivalent to the equality of the local factors:  $L_p(\alpha, s) = L_p(M, s)$  for all *p*, and  $L_{\infty}(\alpha, s) = L_{\infty}(M, s)$ .
- (6) In part (ii), M need not be absolutely irreducible; for example,  $\alpha \in H_1(L_2(N), \mathbb{C})$  could be associated to an elliptic curve with complex multiplication.

Let us return to our example of  $L_1(N)$ , and show how we can use the Conjecture to deduce quadratic reciprocity. The group  $H_0(L_1(N), \mathbb{C})$  is isomorphic to  $H_0((\mathbb{Z}/N\mathbb{Z})^*, \mathbb{C})$ , and this homology group consists of cycles  $\zeta_f = \sum f(x)x$  for functions  $f : (\mathbb{Z}/N\mathbb{Z})^* \to \mathbb{C}$ . Let us work out how the operator  $T_p^{(1)}$  behaves.

Because we are free to scale our lattice by a scalar, we can take  $\Lambda$  to be  $\mathbb{Z}$ , with distinguished generator g = 1 (although, for clarity, we shall continue to write g). The level N structure  $\phi$  is determined by  $\phi(g) = a \in (\mathbb{Z}/N\mathbb{Z})^*$ .

We know that  $T_p^{(1)}(\Lambda, \phi)$  must be the pair  $(p\Lambda, \phi|_{p\Lambda}) = (\Lambda', \psi)$ . Since  $\Lambda'$  has distinguished generator pg, we have  $\psi(pg) = p\psi(g) = pa$ , so

$$T_p^{(1)}\zeta_f = \sum f(x)px = \sum f(p^{-1}x)x,$$

where  $p^{-1}$  is the inverse of p modulo N. Therefore  $(T_p^{(1)}f)(x) = f(p^{-1}x)$ .

Suppose that  $\alpha$  is a simultaneous eigenelement for  $T_p^{(1)}$  for all p not dividing N. We may scale  $\alpha$  so that  $\alpha(1) = 1$ . We know that  $(T_p^{(1)}\alpha)(x) = \alpha(p^{-1}x)$ , and we also have  $(T_p^{(1)}\alpha)(x) = a_p\alpha(x)$  (we need not write  $a_p^{(1)}$ , since this example is onedimensional). Therefore  $\alpha(p^{-1}x) = a_p\alpha(x)$ . Set x = 1, and we have  $\alpha(p^{-1}) = a_p$ . Now set x = p, and we have  $\alpha(p) = a_p^{-1}$ . We can now use induction to conclude that  $\alpha(p^k) = a_p^{-k}$ . This formula holds for the image of any prime p in  $(\mathbb{Z}/N\mathbb{Z})^*$ , which implies that  $\alpha(xy) = \alpha(x)\alpha(y)$ . In other words, the eigenelement  $\alpha$  is a character

of  $(\mathbb{Z}/N\mathbb{Z})^*$ , and the eigenvalue  $a_p$  is  $\alpha(p^{-1})$ . Notice that the eigenfunction  $\alpha$  also determines the eigenvalue  $a_p$ .

Let us move now to the other side of the correspondence. Take any non-zero integer W, and let V be the variety defined by the equation  $x^2 - W = 0$ . We compute the Hasse-Weil L-function  $Z_p(V,t)$  for primes  $p \nmid W$ . If the quadratic residue symbol  $\left(\frac{W}{n}\right) = 1$ , then  $\#V(\mathbf{F}_{p^m}) = 2$  for all positive integers *m*. Therefore

$$Z_p(V,t) = \exp\left(\sum_{m=1}^{\infty} \frac{2}{m} t^m\right)$$
$$= \exp\left(2(-\log(1-t))\right)$$
$$= \frac{1}{(1-t)^2}$$
$$= \frac{1}{(1-t)\left(1-\left(\frac{W}{p}\right)t\right)}.$$

If  $\left(\frac{W}{p}\right) = -1$ , then  $\#V(\mathbf{F}_{p^{2m}}) = 2$  and  $\#V(\mathbf{F}_{p^{2m+1}}) = 0$ . Therefore

$$\begin{split} Z_p(V,t) &= \exp\left(\sum_{m=1}^\infty \frac{2}{2m}t^{2m}\right) \\ &= \exp\left(\sum_{m=1}^\infty \frac{(t^2)^m}{m}\right) \\ &= \exp\left(-\log(1-t^2)\right) \\ &= \frac{1}{1-t^2} \\ &= \frac{1}{(1-t)\left(1-\left(\frac{W}{p}\right)t\right)}. \end{split}$$

We now have

$$\prod_{p} Z_p(V, p^{-s}) = \zeta(s) L(\chi, s),$$

where  $L(\chi, s)$  is the Dirichlet *L*-series defined by the function  $\chi(p) = \left(\frac{W}{p}\right)$ . The 'interesting' part of the cohomology of *V* corresponds to the *L*-series  $L(\chi, s)$ , and this is the finite part of the L-series of the motive M that we study. Precisely, for each p of good reduction, we must make a choice  $S_p$  of a subset of  $\{1, (\frac{W}{p})\}$ , and we select  $S_p = \left\{ \left(\frac{W}{p}\right) \right\}$ . The conductor N of the variety M is a divisor of 4W. (In fact, it might be considerably less than 4|W|, since we have not as yet asked that W be square-free; in addition, if  $W \equiv 1 \pmod{4}$ , then the factor 4 is unnecessary.) The Conjecture tells us that there is an eigenelement  $\alpha \in H_0((\mathbb{Z}/4W\mathbb{Z})^*,\mathbb{C})$  such that  $\alpha(p) = \left(\frac{W}{p}\right)$ . (This is where the fact that the eigenelement  $\alpha$  also determines the eigenvalue  $a_p$  comes into play, along with the helpful observation that we can ignore the exponent of -1 in our formulas for  $\alpha(p)$  because  $-1^{-1} = -1$ . We have also used the fifth observation following the definition of  $T_p^{(k)}$ , with N' = 4|W|.)

Suppose initially that W = -1. Since  $\alpha$  can be thought of as a character on  $(\mathbb{Z}/4\mathbb{Z})^*$ , we can conclude that  $\left(\frac{-1}{p}\right)$  is defined by the residue class of  $p \pmod{4}$ . Since  $\left(\frac{-1}{3}\right) = -1$  and  $\left(\frac{-1}{5}\right) = 1$ , we have deduced the usual formula for  $\left(\frac{-1}{p}\right)$ .

#### AVNER ASH AND ROBERT GROSS

Next take W = 2, and we now have that  $\alpha$  is defined (mod 8). Computation of  $\left(\frac{2}{n}\right)$  for p = 3, 5, 7 and 17 gives the usual formula for  $\left(\frac{2}{n}\right)$ .

Next suppose that  $p \equiv q \pmod{4}$ , with p > q, and let W = (p-q)/4. Then  $p \equiv q \pmod{4W}$ , which means that  $\alpha(p) = \alpha(q)$ , or  $\left(\frac{W}{p}\right) = \left(\frac{W}{q}\right)$ . We have

$$\binom{p}{q} = \binom{4W+q}{q} = \binom{4W}{q} = \binom{W}{q} = \binom{W}{p} = \binom{4W}{p} = \binom{p-q}{p} = \binom{-q}{p} = \binom{-1}{p}\binom{q}{p},$$

which implies most of the usual formula for quadratic reciprocity.

To derive the remaining case, observe that the congruence

$$x^2 - W \equiv 0 \pmod{4W - 1}$$

always has the solution  $x \equiv 2W$ . This tells us that if W is positive and p is any prime dividing 4W - 1, then  $\alpha(p) = 1$ , so  $\alpha(4W - 1) = 1$ .

Now suppose that  $p + q \equiv 0 \pmod{4}$ , and let W = (p + q)/4. Our preceding observation implies that  $\alpha(p) = \alpha(q)$ , which in turn implies that  $\left(\frac{W}{p}\right) = \left(\frac{W}{q}\right)$ , and, reasoning as before, we can conclude that  $\left(\frac{p}{q}\right) = \left(\frac{q}{p}\right)$ .

We next take a more complex example, to explain (finally!) what it means for an elliptic curve to be modular. For instance, we shall consider the curve

$$E : y^2 + y = x^3 - x^2$$

with discriminant  $\Delta = -11$  and conductor N = 11. For any elliptic curve V given by a non-singular cubic equation in the plane, we write  $\hat{V}$  for the complete curve, that is,  $V \cup \{\infty\}$ . We count the number of solutions of this equation modulo p for various primes  $p \neq 11$ , and add 1 for the 'point at  $\infty$ ', to obtain  $\#\hat{V}(\mathbf{F}_p)$ .

For any elliptic curve V defined over  $\mathbf{Q}$ , set  $\#\hat{V}(\mathbf{F}_p) = 1 + p - a_p$ . It is plausible, though not obvious, that one might want to write the number of solutions in this way: for roughly half of the p possible values of x, one expects the left-hand side to have a solution, but in those cases, one can expect there to be 2 solutions, since  $y^2 + y$  is a quadratic polynomial. Therefore one expects roughly p solutions to the congruence, plus an additional solution for the point at infinity, and then we can consider  $a_p$  to measure how far the actual number of solutions deviates from the expected number.

We can also motivate the expression  $1 - a_p + p$  by thinking in terms of motives. The elliptic curve V affords the motives  $H^0(\hat{V})$ ,  $H^1(\hat{V})$  and  $H^2(\hat{V})$ . The motives in dimensions 0 and 2 are essentially trivial, and contribute 1 and p, respectively (via the same analysis that gives the formula  $\#\mathbf{P}^1(\mathbf{F}_{p^m}) = 1 + p^m$ ). The number  $a_p$ corresponds to the non-trivial motive in dimension 1.

The statement that an elliptic curve V of conductor N is 'modular' is equivalent to the fact that there is an eigenclass  $\alpha \in H_1(L_2(N), \mathbb{C})$  such that

$$T_p^{(1)}(\alpha) = a_p \alpha, \quad p \nmid N.$$

If you know the usual definition of 'modular elliptic curve', then you can see this as follows. The set  $H_1(L_2(N), \mathbb{C})$  is closely connected to the 'space of modular forms of weight 2 and level N'. In fact, by a theorem of Eichler and Shimura [18, Chapter 7], any cuspidal H-eigenelement  $\alpha \in H_1(L_2(N), \mathbb{C})$  has the same H-eigenvalues as some newform of weight 2 and level N.

One can check that the elliptic curve E given above is modular by finding the corresponding  $\alpha$ , or equivalently the weight 2 modular form of level 11; see, for

instance, [17, 13]. If we set  $q = e^{2\pi i \tau}$  and

$$\eta(\tau) = e^{\pi i \tau/12} \prod_{n=1}^{\infty} (1-q^n),$$

then the modular form corresponding to E is

$$f(\tau) = \eta(\tau)^2 \eta(11\tau)^2$$
  
=  $\sum a_n q^n$   
=  $q - 2q^2 - q^3 + 2q^4 + q^5 + 2q^6 - 2q^7 - 2q^9 - 2q^{10}$   
+  $q^{11} - 2q^{12} + 4q^{13} + 4q^{14} - q^{15} - 4q^{16} - 2q^{17} + 4q^{18} + 2q^{20} + \cdots$ 

It is amusing to check a few of the  $a_p$  by hand. For instance, if p = 2, then we compute the solutions  $\hat{E}(\mathbf{F}_2) = \{(0,0), (0,1), (1,0), (1,1), \infty\}$ . Therefore  $\#\hat{E}(\mathbf{F}_2) = 5 = 1 + 2 - (-2)$ , so  $a_2$  should equal -2. It does: the coefficient of  $q^2$  in  $f(\tau)$  is -2.

The conjectures stated above imply the well-known 'modularity conjecture': *every elliptic curve defined over*  $\mathbf{Q}$  *is modular.* 

Notice that in our first example, using  $L_1(N)$ , we did not include a 'point at infinity' on V because V was already complete, since it consisted simply of 2 geometric points. In general, including points at infinity compactifies the geometric space and simplifies the formulas.

For our final example, we move to a higher dimension.  $H_3(L_3(N), \mathbb{C})$  has been computed for some small values of N in [1, 2, 10, 9], and certain eigenclasses have been experimentally 'related' to Galois representations. Van Geeman and Top have related a few of these classes to varieties in [10]; here is an example.

Consider  $\hat{V}$  to be the variety given by

$$t^{2} = xy(x^{2} - 1)(y^{2} - 1)(x^{2} - y^{2} + 2xy)$$

along with points at  $\infty$ . Then there is an H-eigenclass in  $H_3(L_3(128), \mathbb{C})$  with

$$T_p^{(k)} \alpha = a_p^{(k)} \alpha, \quad k = 0, 1, 2, 3,$$

so that  $\#\hat{V}(\mathbf{F}_{p^m})$ , for  $m \ge 1$ , corresponds in a precise though rather complicated way to the pair  $(a_p^{(1)}, a_p^{(2)})$  for all primes  $p \le 67$ . (Note that  $a_p^{(0)} = a_p^{(3)} = 1$  for all p.) In particular, define

$$N_{p^m}(V) := \#\{(x, y, z) \in \mathbf{F}_{p^m}^3 \mid t^2 = xy(x^2 - 1)(y^2 - 1)(x^2 - y^2 + 2xy)\},\$$

$$N_{p^m}(E) := \#\{(v, w) \in \mathbf{F}_{p^m}^2 \mid w^2 = v(v^2 + 2v - 1)\}.$$

The second quantity corresponds to points at  $\infty$ , and may be considered 'understood', as *E* is a modular elliptic curve. Then

$$\sum_{i=1}^{6} \alpha_i^m = N_{p^m}(V) + 2N_{p^m}(E) - p^{2m} - 2p^m \left(1 + \left(\frac{2}{p}\right)^m\right),$$

where the numbers  $\alpha_i$  are determined by the equation

6

$$\prod_{i=1}^{6} (1 - \alpha_i X) = X^6 - c_1 X^5 + c_2 X^4 - c_3 X^3 + p^2 c_2 X^2 - p^4 c_1 X + p^6$$
  
=  $(X^3 - \chi(p) b_p X^2 + p \overline{b_p} X - \chi(p) p^3) (X^3 - \chi(p) \overline{b_p} X^2 + p b_p X - \chi(p) p^3)$ 

where  $b_p = a_p^{(1)}$ ,  $\overline{b_p} = a_p^{(2)}$  and  $\chi(p) = \left(\frac{-2}{p}\right)$ . The motive in question here is a 6-dimensional piece of  $H^2(\hat{V})$ , which itself is 34-dimensional.

There is as yet no proof of this correspondence for all p, though there is no reason other than lack of computer time that the computations cannot be continued for larger primes.

The status of the generalized reciprocity conjecture, restricted as we have treated it, that is, for homology with trivial coefficients, is summarized neatly by the number n. For n = 1, it is essentially equivalent to class field theory for  $\mathbf{Q}$  or the theory of cyclotomic fields. For n = 2, the results of Breuil, Conrad, Diamond, Taylor and Wiles cited in the introduction, along with results of Eichler and Shimura [8, 18], give a good piece of the picture. For  $n \ge 3$ , the conjecture is mostly unproven.

## Appendix

In this appendix, we introduce Galois representations, which are really at the heart of the conjectures described above. We also give one more example; in this one, the motive conjectured to exist is not yet known.

A motive affords a compatible series of  $\ell$ -adic representations of  $G_Q$ , the Galois group of  $\overline{Q}$  over Q. In an  $\ell$ -adic representation  $\rho$ , for each finite unramified prime p, the characteristic polynomial  $F_p$  of a Frobenius element  $\phi_p$  at p is well-defined and gives information about how  $\rho(\phi_p)$  acts in the representation. Thus a formula expressing  $F_p$  in some other terms can be thought of as a generalized reciprocity law. Moreover, as we have seen above,  $F_p$  is closely related to the  $\#V(\mathbf{F}_{p^m})$  for the variety V from which the given motive comes.

For a last concrete example, consider the Hecke eigenclass  $\alpha$  on  $L_3(61)$  from [1]. If we set  $\omega = (1 + \sqrt{-3})/2$ , then the first few  $a_p^{(1)}$  were computed to be as follows.

р	2	3	5	7	11
$a_p^{(1)}$	$1-2\omega$	$-5+4\omega$	$-2+4\omega$	$-6\omega$	$-2+2\omega$

For each p,  $a_p^{(2)}$  is the complex conjugate of  $a_p^{(1)}$ . We then conjecture the existence of a continuous 3-dimensional  $\lambda$ -adic representation  $\rho$  of  $G_Q$  (where  $\ell$  is a rational prime, and  $\lambda$  is a prime in some number field above  $\ell$ ) unramified outside 61 $\ell$  such that

$$F_p = X^3 - a_p^{(1)}X^2 + pa_p^{(2)}X - p^3.$$
(\*)

In [3], we reduced all the coefficients modulo  $\sqrt{-3}$ , and looked for a  $\overline{\rho} : \mathbf{G}_{\mathbf{Q}} \to \mathbf{GL}(3, \mathbf{F}_3)$ . This would be the weakest possible check of the conjecture for  $\ell = 3$ . Knowing the  $\overline{F_p}$  allowed us to infer how  $\overline{\rho}(\phi_p)$  must act, and hence (eventually) how p must split in the fixed field of ker  $\overline{\rho}$ . For instance, if p = 11, then  $a_p^{(1)} \equiv a_p^{(2)} \equiv -1 \pmod{3}$  and  $p \equiv -1 \pmod{3}$ , so  $F_p \equiv X^3 + X^2 + X + 1$ . Thus  $\overline{\rho}(\phi_p)$  must be a matrix in  $\mathbf{GL}(3, \mathbf{F}_3)$  with that characteristic polynomial. Juggling such information, we determined the only possible  $\overline{\rho}$ , and showed that it matched the given data. We have no idea how to show that  $\overline{\rho}$  is really attached to  $\alpha$  (that is, that (\*) holds for all  $p \nmid 3 \cdot 61$ ), nor how to find  $\rho$ .

To summarize this example: we have here reciprocity between a certain series of representations of  $G_Q$  and the Hecke information contained in  $\alpha$ . It is relatively easy to compute  $\alpha$  and then to conjecture these rather deep properties of  $G_Q$ .

In the case of the proof of Fermat's last theorem, a putative solution of the Fermat equation leads to a certain representation of  $G_Q$ . It is this representation

that is shown not to exist by proving that it should be controlled through reciprocity by a certain Hecke eigenclass in  $L_2(2)$ , which is readily checked to be non-existent.

ACKNOWLEDGEMENTS. We should like to thank Fred Diamond, Nicholas Katz, Barry Mazur, David Rohrlich, Joseph Silverman, Glenn Stevens and the referee for comments on earlier drafts of this paper.

#### References

- 1. A. ASH, D. GRAYSON and P. GREEN, 'Computations of cuspidal cohomology of congruence subgroups of SL(3, Z)', J. Number Theory 19 (1984) 412-436.
- 2. A. ASH and M. MCCONNELL, 'Double cuspidal cohomology for principal congruence subgroups of GL(3, Z)', Math. Comp. 59 (1992) 673-688.
- 3. A. ASH, R. PINCH and R. TAYLOR, 'An  $\hat{A}_4$  extension of Q attached to a non-selfdual automorphic form on GL(3)', Math. Ann. 291 (1991) 753-766.
- 4. L. CLOZEL, 'Motifs et formes automorphes: applications du principe de fonctorialité', Automorphic Forms, Shimura Varieties, and L-functions, Proc. Ann Arbor Conf. (ed. L. Clozel and J. S. Milne, Academic Press, New York, 1990) 77-159.
- 5. B. CONRAD, F. DIAMOND and R. TAYLOR, 'Modularity of certain potentially Barsotti-Tate Galois representations', J. Amer. Math. Soc. 12 (1999) 521-568.
- 6. H. DARMON, 'A proof of the Full Shimura-Taniyama-Weil Conjecture is announced', Notices Amer. Math. Soc. 46 (1999) 1397-1401.
- F. DIAMOND, 'On deformation rings and Hecke rings', Ann. of Math. 144 (1996) 137–166.
   M. EICHLER, 'Quaternäre quadratische Formen und die Riemannsche Vermutung für die Kongruenzzetafunktion', Acta Math. 5 (1954) 355-366.
- 9. B. VAN GEEMAN, W. VAN DER KALLEN, J. TOP and A. VERBERKMOES, 'Hecke eigenforms in the cohomology of congruence subgroups of SL(3, Z)', Experiment. Math. 6 (1997) 163-174.
- 10. B. VAN GEEMAN and J. TOP, 'A non-selfdual automorphic representation of GL<sub>3</sub> and a Galois representation', Invent. Math. 117 (1994) 391-401.
- 11. S. GELBART, 'Elliptic curves and automorphic representations', Adv. Math. 21 (1976) 235-292.
- 12. U. JANNSEN, S. KLEIMAN and J.-P. SERRE, Motives, Proc. Sympos. Pure Math. 55, Parts 1 and 2 (Amer. Math. Soc., Providence, RI, 1994).
- 13. A. W. KNAPP, Elliptic curves, Math. Notes 40 (Princeton University Press, Princeton, NJ, 1992).
- 14. R. P. LANGLANDS, 'Some contemporary problems with origins in the Jugendtraum (Hilbert's problem 12)', Hilbert problems, Proc. Sympos. Pure Math. 28, Part 2 (Amer. Math. Soc., Providence, RI. 1976) 401-418
- 15. K. RIBET, 'Galois representations and modular forms', Bull. Amer. Math. Soc. 32 (1995) 375-402.
- 16. J.-P. SERRE, A course in arithmetic, Grad. Texts in Math. 7 (Springer, New York, 1973).
- 17. G. SHIMURA, 'A non-solvable reciprocity law', J. Reine Angew. Math. 221 (1966) 209-220.
- 18. G. SHIMURA, Introduction to the arithmetic theory of automorphic forms (Princeton University Press, Princeton, NJ, 1971).
- 19. J. SILVERMAN, The arithmetic of elliptic curves, Grad. Texts in Math. 106 (Springer, New York, 1986).
- J. TATE, 'Problem 9: the general reciprocity law', *Hilbert problems*, Proc. Sympos. Pure Math. 28, Part 2 (Amer. Math. Soc., Providence, RI, 1976) 311–322. 21. R. TAYLOR and A. WILES, 'Ring theoretic properties of certain Hecke algebras', Ann. of Math. 141
- (1995) 553-572.
- 22. A. WEIL, Number theory: an approach through history from Hammurapi to Legendre (Birkhäuser, Boston, 1984).
- 23. A. WILES, 'Modular elliptic curves and Fermat's Last Theorem', Ann. of Math. 141 (1995) 443-551.
- 24. B. WYMAN, 'What is a reciprocity law?', Amer. Math. Monthly 79 (1972) 571-586.

Ohio State University Columbus, OH 43210-1174 USA

Boston College Chestnut Hill, MA 02467-3806 USA